# GigaGS: Scaling up Planar-Based 3D Gaussians for Large Scene Surface Reconstruction

**Junyi Chen**[1,2,*], **Weicai Ye**[1,3*✉], **Yifan Wang**[1,2], **Danpeng Chen**[3], **Di Huang**[1],
**Wanli Ouyang**[1], **Guofeng Zhang**[3], **Yu Qiao**[1], **Tong He**[1,✉]

[1]Shanghai AI Laboratory [2]Shanghai Jiao Tong University [3]State Key Lab of CAD&CG, Zhejiang University
https://open3dvlab.github.io/GigaGS/

## Abstract

3D Gaussian Splatting (3DGS) has shown promising performance in novel view synthesis. Previous methods adapt it to obtaining surfaces of either individual 3D objects or within limited scenes. In this paper, we make the first attempt to tackle the challenging task of large-scale scene surface reconstruction. This task is particularly difficult due to the high GPU memory consumption, different levels of details for geometric representation, and noticeable inconsistencies in appearance. To this end, we propose GigaGS, the first work for high-quality surface reconstruction for large-scale scenes using 3DGS. GigaGS first applies a partitioning strategy based on the mutual visibility of spatial regions, which effectively grouping cameras for parallel processing. To enhance the quality of the surface, we also propose novel multi-view photometric and geometric consistency constraints based on Level-of-Detail representation. In doing so, our method can reconstruct detailed surface structures. Comprehensive experiments are conducted on various datasets. The consistent improvement demonstrates the superiority of GigaGS.

## 1   Introduction

3D Gaussian Splatting (Kerbl et al. 2023) has demonstrated remarkable performance on the task of novel view synthesis. Recently, some methods (Guédon and Lepetit 2023; Huang et al. 2024a) adapted it to surface reconstruction, which is drawing increasing attention for its numerous promising applications, such as 3D asset generation (Tang et al. 2023; Chen et al. 2024b; He et al. 2024; Xu et al. 2024; Tang et al. 2024a) and virtual reality (Wu et al. 2023; Charatan et al. 2024). Compared with the task of novel view synthesis, recovering the inherited 3D surfaces is much more challenging as it requires preserving 3D coherence throughout varying perspectives with only 2D projects for supervision.

Although significant advances (Huang et al. 2024a; Guédon and Lepetit 2023; Yu, Sattler, and Geiger 2024) have been made, these methods either focus on object-level reconstruction or struggle to capture intricate geometric surfaces. These challenges become even more critical in the context of large-scale scene reconstruction. Firstly, the computational resource consumption is enormous, as a scene covering several square kilometers often contains billions of Gaussian points.

Directly applying previous methods may lead to suboptimal reconstruction quality or run into memory-related issues. Secondly, previous neural reconstruction methods with 3D Gaussian Splatting (Kerbl et al. 2023) often address the challenging task by introducing monocular regularization. For example, SuGaR (Guédon and Lepetit 2023) introduces single-view depth and normal geometry consistency constraints to ensure the correctness of single-view geometry. Despite achieving good reconstruction results, these methods fail to incorporate multi-view constraints to ensure global geometry consistency. Recent work of PGSR (Chen et al. 2024a) first utilizes multiview geometric constraints to 3DGS representation. Although impressive, the method fails to capture the geometric details at different scales.

To address the above challenges, we propose GigaGS. To the best of our knowledge, it is the first work of high-quality surface reconstruction for large-scale scenes using 3DGS. Firstly, we implement an efficient and scalable partitioning strategy to address the computational demands of processing large-scale scenes. Unlike conventional approaches relying on spatial distance metrics, we introduce a novel grouping mechanism based on the mutual visibility of spatial regions captured by the scene cameras. This enables us to partition the scene into overlapping blocks that can be processed in parallel. Each block undergoes independent optimization, thus allowing for distributed processing of the scene data. Subsequently, the optimized blocks are seamlessly merged to reconstruct the complete scene, ensuring computational efficiency without compromising on reconstruction accuracy. Secondly, we present a novel method to harness multi-view photometric and geometric consistency constraints within a Level-of-Detail (LoD) framework. This approach is designed to enhance the preservation of geometric details across different scales of the reconstructed scene. By integrating LoD representation into the constraint formulation, we ensure that the reconstruction process maintains fidelity and coherence across varying levels of scene complexity. Leveraging both photometric and geometric information from multiple views, our method facilitates robust reconstruction of intricate scene details while mitigating artifacts and inconsistencies.

To summarize, the contributions of the paper are listed as follows:

- To the best of our knowledge, we are the first to utilize
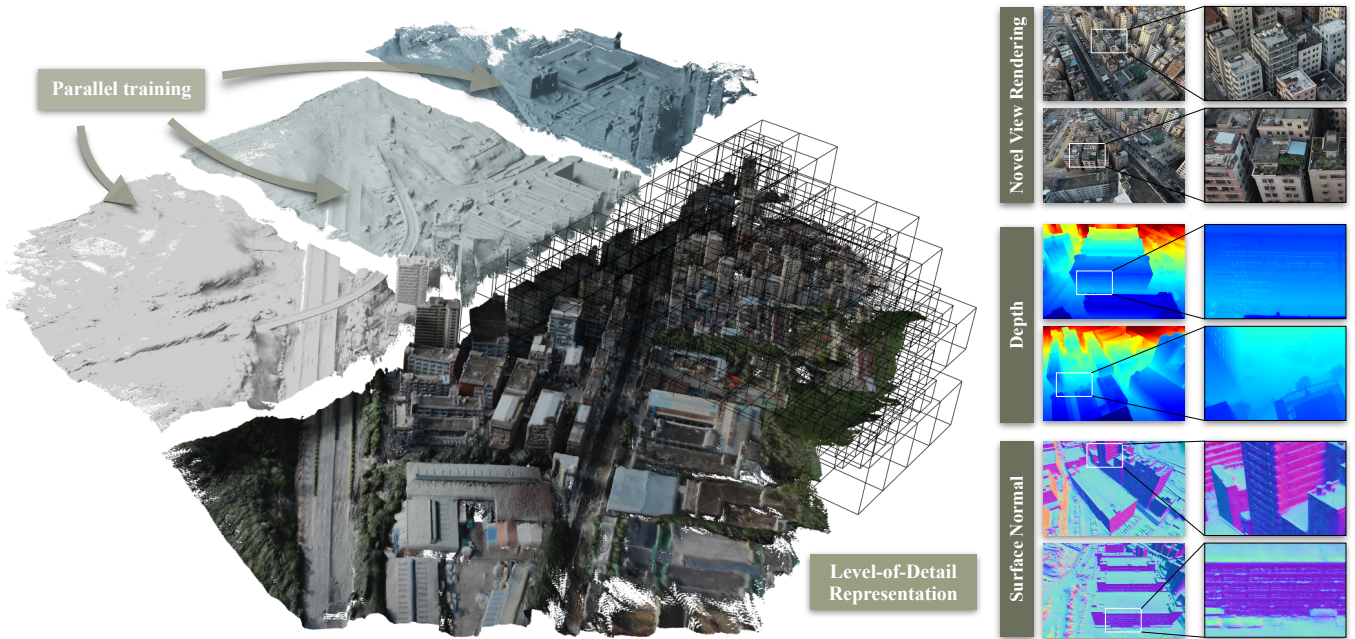
---

*Equal Contribution.

Figure 1: We propose GigaGS, the first work specifically designed for large scene surface reconstruction. Our approach ensures high rendering quality while also extracting high-quality meshes.

3DGS for large-scale surface reconstruction.

- Based on a large-scene partitioning strategy, we present a novel method to add multi-view photometric and geometric consistency constraints within a Level-of-Detail (LoD) framework.

- Comprehensive experiments on various datasets demonstrate the effectiveness of the proposed method for large scene surface reconstruction.

## 2 Related Work

### 2.1 Neural Rendering

Neural radiance field (Mildenhall et al. 2021; Ye et al. 2023a; Huang et al. 2024b; Ming, Ye, and Calway 2022) models a 3D scene by learning a continuous volumetric scene function that maps 3D coordinates and viewing directions to the corresponding RGB color and volume density. This approach enables the synthesis of novel views of complex scenes from a set of input images. Despite the numerous advancements made in recent studies (Fridovich-Keil et al. 2022; Chen et al. 2022; Sun, Sun, and Chen 2022; Müller et al. 2022) aimed at enhancing its performance, such as reducing training duration and expediting rendering processes, the existing framework remains constrained by the absence of a clear and explicit representation. Consequently, this limitation poses significant challenges in extending its applicability to a broader spectrum of scenarios. 3D Gaussian Splatting (Kerbl et al. 2023) is another innovative approach in neural rendering that uses Gaussian functions to represent volumetric data. This technique involves placing 3D Gaussians in the scene to approximate the spatial distribution of radiance and density. More importantly, 3DGS

possesses an explicit representation that empowers real-time rendering by utilizing rasterized rendering methods. Subsequent works focus on enhancing rendering quality(Yu et al. 2023; Lu et al. 2023), further streamlining 3DGS to improve rendering speed(Fan et al. 2023), and extending its application to reflective surfaces(Jiang et al. 2023).

### 2.2 Surface Reconstruction

Surface reconstruction (Chen et al. 2024a; Ye et al. 2024b,a; Tang et al. 2024b; Ye et al. 2022, 2023b; Liu et al. 2021; Li et al. 2020) aims at generating accurate and detailed 3D models from various forms of input data, which is fundamental for numerous applications, including 3D modeling, virtual reality, and robotics. Recent advancements in deep learning have significantly improved the quality and efficiency of surface reconstruction techniques. The utilization of neural implicit representation offers the advantage of continuous and differentiable surfaces, thereby enabling more precise and flexible reconstruction. Recent research (Li et al. 2023; Guo et al. 2022; Wang et al. 2021; Yu et al. 2022) has extensively employed these representations, demonstrating their ability to generate detailed and accurate surface reconstructions capable of capturing intricate geometric details and complex topological structures. SuGaR (Guédon and Lepetit 2023) successfully aligned 3DGS with the surface of the scene, yielding remarkable reconstruction outcomes. Additionally, 2DGS (Huang et al. 2024a) simplifies 3DGS, thereby facilitating a more expressive representation of the scene's structure. However, there remains scope for further enhancement in achieving quantitative results.

## 2.3 large scale reconstruction

Large scale reconstruction involves creating detailed and accurate 3D models of extensive environments, and is challenging due to the vast amount of data, the need for high precision, and the complexity of the scenes. MegaNeRF (Turki, Ramanan, and Satyanarayanan 2022) represents a pioneering approach for reconstructing expansive outdoor scenes. It extends the NeRF framework by partitioning the scene into manageable blocks and independently optimizing each block, thus enabling effective handling of large-scale environments. Furthermore, VastGaussian (Lin et al. 2024) proposed a blocking strategy specifically tailored for 3DGS, enabling parallel training of distinct blocks. This strategy not only reduces training time but also facilitates the attainment of high-quality rendering for the entire scene.

## 3 Preliminaries

In this work, we employ 3DGS (Kerbl et al. 2023) as the fundamental 3D representation and rendering entity, while employing unbiased depth rendering to acquire depth maps and surface normal maps. Within this section, we shall elucidate the significance of these two technologies as essential contextual foundations.

### 3.1 3D Gaussian Splatting

One of the core parts in the 3DGS is the 3D Gaussian kernel, which encapsulates the visual characteristics of a spatial region. Each 3D Gaussian possesses several key attributes, including positional coordinates, a covariance matrix that describes the arrangement of the kernel, opacity, and spherical harmonic coefficients that encode the view-dependent colors. During the rendering procedure, the 3D Gaussians are projected onto a 2D Gaussian distribution specific to the given viewpoint. Subsequently, the final rendering output for that viewpoint is generated through $\alpha$-blending:

$$C = \sum_{i \in M} c_i \alpha_i T_i, \quad T_i = \prod_{j=1}^{i-1}(1 - \alpha_j), \qquad (1)$$

The parameters are updated via a differentiable rendering process.

### 3.2 Unbiased Depth Rendering

2DGS (Huang et al. 2024a) represents 3D shape with 2D Gaussian primitives. SuGaR (Guédon and Lepetit 2023) adds a regularization term that encourages the Gaussians to align with the surface of the scene. It inherently provides accurate estimations of the surface normal, which corresponds to the shortest axis of the Gaussian kernel. Inspired by these methods, PGSR (Chen et al. 2024a) was initially developed for the purpose of viewpoint-dependent normal vector rendering:

$$N = \sum_{i \in M} R_c n_i \alpha_i T_i, \qquad (2)$$

where $R_c$ is the rotation matrix from camera coordinates to world coordinates and $n_i$ is the normal vector of i-th 3DGS. Unlike previous approaches, PGSR takes a different approach by not directly rendering based on the spatial position of the 3DGS kernels. Instead, it assumes that the 3D Gaussian kernels can be flattened into a plane and fitted onto the actual surface. It then proceeds to render the distance from the camera origin to this Gaussian plane, denoted as $\mathscr{D}$:

$$\mathscr{D} = \sum_{i \in M} d_i \alpha_i T_i, \qquad (3)$$

where, $d_i = (R_c^T(\mu_i - T_c))R_c^T n_i^T$ represent the distance from the camera origin to $i$-th Gaussian Kernel. Once the distances and normals of the planes are obtained, PGSR determine the corresponding depth map by intersecting rays with these planes. This intersection operation ensures that the depth shapes align with the planes assumed by the Gaussian kernel, resulting in a depth map that accurately reflects the actual surfaces:

$$D(p) = \frac{\mathscr{D}}{N(p)K^{-1}\widetilde{p}}, \qquad (4)$$

where $p$ is the 2D position in the image plane. $\widetilde{p}$ denotes the uniform coordinates of $p$ and $K$ is the intrinsic coordinates of the camera.

## 4 Method

The primary challenges in large scene surface reconstruction tasks are the vast area of the scene, an excessive number of fine details, and the drastic fluctuation in image brightness caused by changes in lighting and exposure factors. Existing surface reconstruction methods (Li et al. 2023; Wang et al. 2021; Guédon and Lepetit 2023; Guo et al. 2022) primarily focus on small-scale and object-centric scenes and lack explicit designs to address the challenges posed by scaling up, resulting in limited applicability to large-scale scenes. Some works (Lin et al. 2024; Turki, Ramanan, and Satyanarayanan 2022; Li et al. 2024) design complicated strategies for data partitioning, aiming to reduce training time and memory burden. However, these works mainly concentrate on the image rendering quality while disregarding the scene surface.

To address the scalability issue in surface reconstruction tasks, we present an efficient and scalable scene partitioning strategy for parallel training of different partitions across multiple GPUs. To capture fine-grained details across multiple levels of granularity, our framework employs a hierarchical plane representation to store different levels of details (LoD) and achieve high-quality surface reconstruction.

Generally speaking, in Section 4.1, we present our hierarchical plane representation. Subsequently, in Section 4.2, we elaborate on our scalable partitioning strategy, which is based on our representation and circumvents the limitations imposed by hardware and training time, allowing us to utilize a larger number of 3D Gaussian to represent the large scene, even reaching the giga-level scale. Furthermore, we detail how to accurately fit the surface of the scenes in section 4.3 and subsequently extract the mesh in section 4.4.

## 4.1 Hierachical Plane Representation

In typical 3D reconstruction tasks, training data often includes information about the same object at different scales, especially in aerial images. Existing works (Guédon and Lepetit 2023; Wang et al. 2021; Li et al. 2023; Chen et al. 2024a; Yu et al. 2024; Yariv et al. 2021) struggle to directly capture features at different scales because they lack explicitly designed structures to capture levels of details. Therefore, we introduce a new representation combining a hierarchical structure to model the surface of the scene and a flatten form closing to a planar surface.



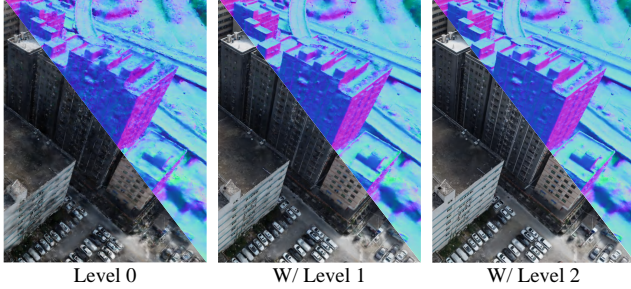| Level 0 | W/ Level 1 | W/ Level 2 |

Figure 2: **Visualization of the effects at different levels.** We visualized the rendered images and normal maps obtained by rendering different levels of the same scene as rendering entities.

**Hierachical Scene Structure**   Given the inherent difficulties in achieving real-time rendering across various scales, especially for conventional 3DGS methods, the surface reconstruction of large-scale scenes has become a challenging task. We adopt a hierarchical structure, inspired by OctreeGS (Ren et al. 2024; Lu et al. 2023), to represent the surface of the scene, where a local set of $n$ 3D Gaussians is represented by a single anchor Gaussian, and during forward inference, the Multi-Layer Perceptron (MLP) is used to recover the parameters of these $n$ 3D Gaussians. The parameters of the MLP are trained jointly with the features of the anchor Gaussians. Different levels of anchor Gaussians are employed to represent features at various levels of granularity. Prior to training, it is feasible to construct a collection of anchor Gaussians at distinct levels from the point cloud $\mathbb{P}$ obtained through Structure-from-Motion (SFM) (Schonberger and Frahm 2016):

$$\text{level}_i = \left\{ v_i \left\lceil \frac{\boldsymbol{p}}{v_i} \right\rceil \Big| \, \boldsymbol{p} \in \mathbb{P} \right\}, \ \ 0 \leq i \leq K-1, \quad (5)$$

where $v_0$ is the fundamental voxel size, and $v_i = v_0/k^i$ is the voxel size of level $i$ and $k$ is the fork number. When $k = 2$, the hierarchical strategy yields the octree structure. $K$ is the maximum number of levels. During the rendering stage, the visibility of levels is determined based on their positional distance $d$ from the viewpoint:

$$\text{upper\_level}(d) = \max(\lceil \log d_{max} - \log d \rceil, K-1). \quad (6)$$

The parameter $d_{max}$ represents the maximum distance between points in the point cloud $\mathbb{P}$ and can be calculated prior to training, and $\lceil \cdot \rceil$ is used to denote the rounding operation. With the distance $d$ increases, fewer 3D Gaussians with lower level are involved in the rendering process. During the training process, we follow the approach of OctreeGS (Ren et al. 2024) for the addition and removal operations of anchor Gaussians.

**Flattened 3D Gaussians**   As the shortest axis of 3D Gaussian kernel inherently provides accurate estimations of the normal vectors (Guédon and Lepetit 2023), we make efforts to compress the minimum axis of each Gaussian kernel during the training phase:

$$\mathcal{L}_{flatten} = \frac{1}{|\mathbb{M}|} \sum_{i \in \mathbb{M}} \big| \min(\boldsymbol{s}_i) \big|, \quad (7)$$

where $\mathbb{M}$ is the set containing the 3D Gaussians. By doing so, we aim to constrain the shortest axis of the 3d Gaussian kernel to be perpendicular to the surface of the scene, thereby facilitating the use of 3DGS to fit the surface of the scene.

## 4.2 Partitioning Strategy

To address the challenges of scaling 3DGS to large scenes, VastGaussian (Lin et al. 2024) proposed a partitioning strategy to evenly distribute the training workload across multiple GPUs. These strategies aim to overcome the limitations of 3DGS in handling large-scale scenes. However, it still face difficulties in extreme scenarios where supervision of a specific region exists in other partitions but is not included in the training data of the current partition due to threshold-based selection strategies. This is likely to occur because aerial trajectories change with the scene, leading to the aforementioned situation.

To tackle this issue, we propose a more robust partitioning approach that leverages the octree-based scene representation. Firstly, we ensure an approximately equal number of cameras in each partition $c$ by following the filtering rule of uniform camera density (Lin et al. 2024). Additionally, each partition is non-overlapping. We have eliminated the manual threshold setting in visibility-based camera selection, which was proposed in VastGaussian Instead, we utilize the painter algorithm (Newell, Newell, and Sancha 1972) to select the cameras based on the partition anchors that can successfully project onto the camera's image plane. This approach aims to cover as many cameras as possible during the selection process. To ensure sufficient supervision for each partition, we project the anchors of the partition onto the image planes of all cameras and consider the camera that can observe the anchors of the partition to be included in the training of that partition. This process is a greedy one without manually setting any thresholds, but it guarantees maximum supervision for each partition.

Finally, we need to expand the partitions so that each camera can render a complete image. Therefore, we project all anchors onto the image plane of each camera and add all visible anchors to the training set of that partition based on the

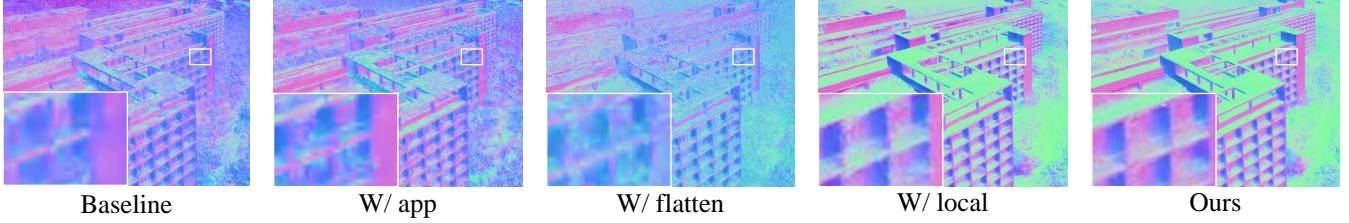| Baseline | W/ app | W/ flatten | W/ local | Ours |

Figure 3: **Different loss terms on the final optimization process.** Initially, using only the image loss as supervision does not yield a satisfactory surface reconstruction. The incorporation of an appearance model reduces certain artifacts. However, the addition of flatten regularization on the geometric structure without additional geometric supervision leads to a decrease in expressiveness by the model. Nevertheless, the inclusion of the local loss allows for improved surface quality. Finally, with the introduction of multi-view regularization, the surface reconstruction performance is further enhanced, highlighting the superiority of our method in surface reconstruction.

equation 6. Since we only require the anchors within the partitions for our final purpose, the expanded anchors are used solely to assist in training. Therefore, in Equation 6, we employ a floor operation rather than a round operation to reduce the number of anchors outside the partitions. This operation does not affect the final quality because those anchors will be discarded after training is completed.

### 4.3 Appearance and Geometry Regularization

As discussed in NeRF-W (Martin-Brualla et al. 2021), directly applying existing representations to a collection of outdoor photographs can lead to inaccurate reconstruction due to factors such as exposure and lighting conditions. These reconstructions exhibit severe ghosting, excessive smoothness, and further artifacts. Therefore, we introduce an appearance model to capture the variations in appearance for each image. Similarly, this model is learned jointly with our planar representation. Next, to ensure that the flattened 3D Gaussians adhere to the actual surface, we enforce consistency between the unbiased depth maps and normal maps from each viewpoint. Simultaneously, we discover that explicit control of consistency across viewpoints has a positive impact on the final surface reconstruction quality.

**Appearance Modeling**   We learn to model the appearance variations in each image in a low-dimensional latent space (Martin-Brualla et al. 2021), such as exposure, lighting, weather, and post-processing effects. In this approach, we allocate an embedding $emb_v$ for each training perspective $v$ and additionally train an appearance model $\phi$ that maps $e$ to per-pixel color adjustment values for the image. By multiplying these adjustment values with the rendered image $\boldsymbol{I}$, we obtain a simulated lighting image that accurately represents the lighting conditions for that particular perspective:

$$\boldsymbol{I}_a = \phi(\boldsymbol{I}, emb_v)\boldsymbol{I}. \tag{8}$$

This enables us to effectively account for the variations in lighting and appearance across different viewpoints, and the following loss is utilized:

$$\mathcal{L}_{app} = L_1(\boldsymbol{I}_a, \boldsymbol{I}_0) + \lambda SSIM(\boldsymbol{I}, \boldsymbol{I}_0), \tag{9}$$

where $\boldsymbol{I}_0$ is the ground-truth image. $\lambda$ is utilized to adjust the relative weights between the two components. Through-

out all our experiments, the value of $\lambda$ is consistently maintained at $0.25$.

**Geometry Consistency**   The vanilla 3DGS (Kerbl et al. 2023), which primarily relies on image reconstruction loss, tends to encounter challenges in local overfitting optimization. In the absence of effective regularization techniques, the intrinsic capacity of 3DGS to capture visual appearance gives rise to inherent ambiguity in the relationship between three-dimensional shape and brightness. Consequently, this ambiguity permits the acceptance of degenerate solutions, resulting in a discrepancy between Gaussian shape estimation and the actual surface representation (Zhang et al. 2020). To address this issue, we have introduced a straightforward regularization constraint, aimed at enforcing geometric consistency between local depth map and surface normal map:

$$\mathcal{L}_{local} = \frac{1}{|\boldsymbol{I}|} \sum_{i \in \boldsymbol{I}} \left| \frac{(P_{i,0} - P_{i,1}) \times (P_{i,2} - P_{i,3})}{|(P_{i,0} - P_{i,1}) \times (P_{i,2} - P_{i,3})|} - \boldsymbol{n}_i \right| \omega_i,$$
$$\tag{10}$$

where $P_{i,j}$ is the 3D position of adjacent pixel $j$ at the top, bottom, left, and right positions relative to pixel $i$ in the camera coordinate system. Additionally, $n_i$ is the normal value of pixel $i$. This assumption holds notable significance as it asserts the interconnectedness of adjacent pixels within a common plane. For pixels that belong to depth discontinuities, we introduce an uncertainty factor, denoted by $\omega_i$, which serves to quantify the likelihood that the pixel $i$ belongs to the boundary of the surface:

$$\omega_i = \left| (P_{i,0} - P_{i,1})(P_{i,2} - P_{i,3}) \right|. \tag{11}$$

Evidently, the dot product mentioned above will yield diminished values in regions characterized by substantial depth disparities, thereby identifying them as edge regions. Consequently, these areas should experience a reduced influence from the constraint imposed by the depth and normal consistency.

**Multi-View Consistency**   Single-view geometry regularization can maintain consistency between depth and normal geometry, but the geometric structures across multiple views are not entirely consistent as shown in figure 3. Therefore,
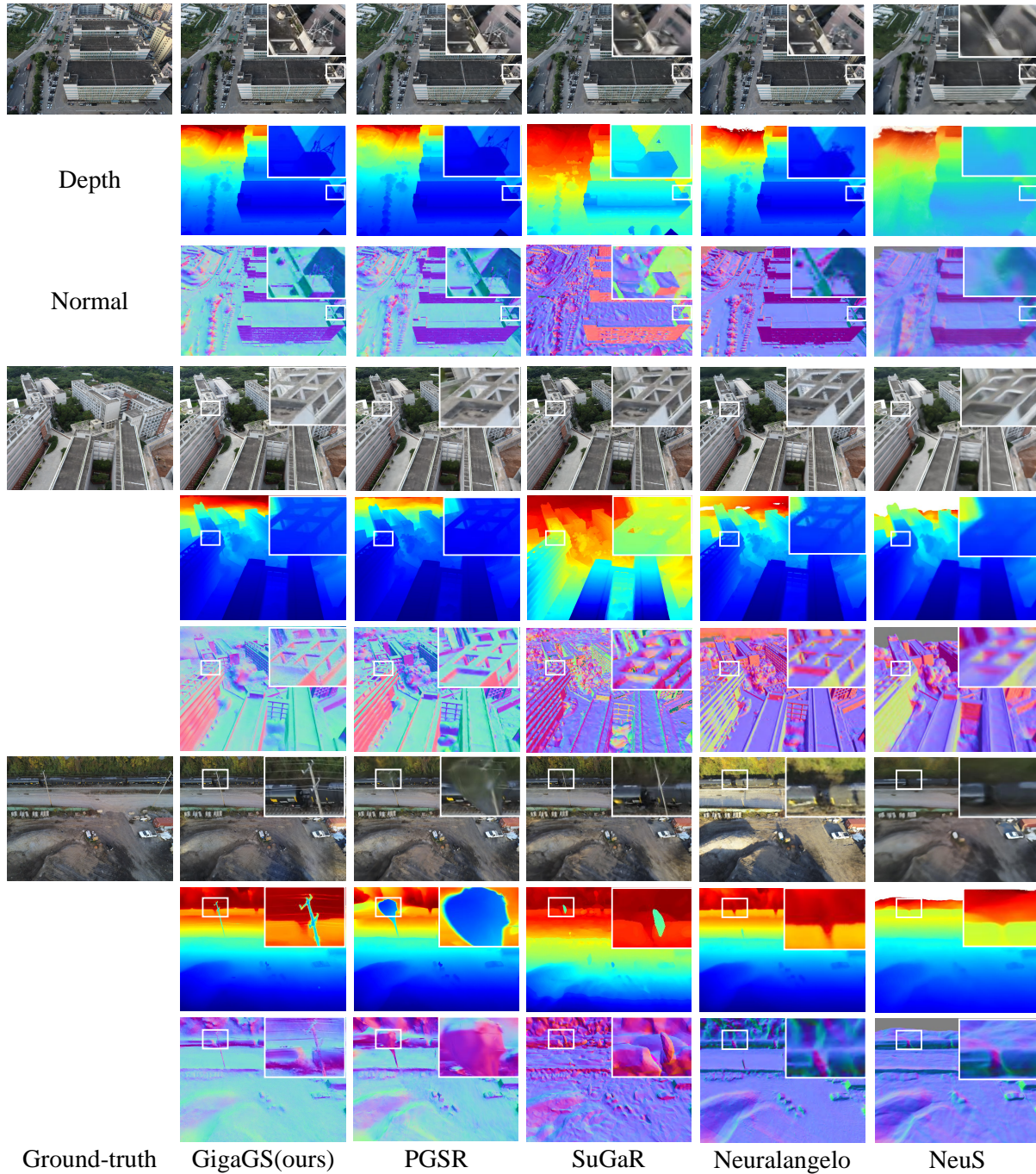
Figure 4: **Comparison of visualization results.** We presented rendered views of the test scenes, along with the corresponding depth maps and normal maps from the same viewpoint.

it is necessary to introduce multi-view geometry regularization to ensure global consistency of the geometric structure. We employ a photometric multi-view consistency constraint based on planar patches to supervise the geometric structure. Specifically, we can render the normal and the distance from each pixel to the plane. Then, the optimization of these geometric parameters can be achieved through patch-based inter-view geometry consistency. For each pixel point $p_r$, we can warp it to an adjacent viewpoint:

$$\tilde{p}_n = H_{rn}\tilde{p}_r, \quad (12)$$

where $\tilde{p}$ is the homogeneous coordinate of pixel point $p$, and homography $H_{rn}$ can be computed as:

$$H_{rn} = K_r(R_{rn} - \frac{T_{rn}n_r^T}{d_r})K_r^{-1}, \quad (13)$$

where $R_{rn}$ and $T_{rn}$ are the relative transformation from the reference frame to the neighboring frame. With a focus on geometric details, we convert the color image $I$ to grayscale image $I$ to supervise our geometric parameters. Then, we utilize the normalized cross correlation (NCC) (Yoo and Han 2009) of patches in the reference frame and the neighboring frame as a metric to evaluate the photometric consistency.

$$\mathcal{L}_{ncc} = \frac{1}{|\mathbb{V}|}\sum_{p_r \in \mathbb{V}}\Big(1 - NCC\big(I(p_r), I(H_{nr}p_n)\big)\Big). \quad (14)$$

Where $\mathbb{V}$ is the valid region checked through geometric consistency constraints. The warp operation may introduce inconsistencies due to occlusions. Therefore, we re-warp the neighboring frames back to the reference frame, and utilize a threshold to filter out areas with significant errors. These areas, where the errors exceed the threshold, are considered as occluded regions:

$$\mathcal{L}_{geo} = \frac{1}{|\mathbb{V}|}\sum_{p_r \in \mathbb{V}}||\tilde{p}_r - H_{nr}H_{rn}\tilde{p}_r||. \quad (15)$$

Finally, the multi-view consistent constrain consists of two components, the multi-view photometric constraint and the multi-view geometric consistency constraint:

$$\mathcal{L}_{mv} = \mathcal{L}_{ncc} + \mathcal{L}_{geo}. \quad (16)$$

As shown in Figure 3, our method demonstrates that multi-view regularization is crucial for reconstruction accuracy. In summary, our overall set of constraints is as follows:

$$\mathcal{L} = \mathcal{L}_{flatten} + \mathcal{L}_{app} + \mathcal{L}_{local} + \mathcal{L}_{mv}. \quad (17)$$

### 4.4 Mesh Extraction

By incorporating our regularization term, we facilitate the generation of a mesh from the optimized Gaussian model. Subsequently, we proceed to render both a visual rendering and a depth map from various vantage points. These rendered images and depth maps are then utilized to fusion into a projected truncated signed distance function (TSDF) volume (Zeng et al. 2017) to finally create the superior quality 3D surface meshes and point clouds.

## 5 Experiments

**Implement Details** In our experiment, we reduced the side length of 4K aerial images to one-fourth of their original size and aligned them with a comparative method. Subsequently, we employed pixel-sfm (Lindenberger et al. 2021) to obtain an initial point cloud from the aerial images and performed Manhattan world alignment, aligning the $y$-axis perpendicular to the world coordinate axis of the ground plane. We divided the entire scene into $4 \times 2$ partitions in the case of rubble, building, residence, and sci-art, while for the largest scene, campus, we divided it into $4 \times 4$ partitions. Each partition was subjected to training for 120,000 iterations to ensure sufficient convergence. Upon the completion of independent training for each partition, we discard all anchors except for those in the original partition $c$. This approach ensures that each partition is ultimately non-overlapping, thereby enabling the construction of a comprehensive scene.

**Baseline Methods** For the purpose of comparing our surface reconstruction results, we have selected Neuralangelo (Li et al. 2023), NeuS (Wang et al. 2021), PGSR (Chen et al. 2024a), and SuGaR (Guédon and Lepetit 2023) as the comparative methods. Neuralangelo and Neus are methods rooted in the Nerf framework, whereas Sugar and PGSR are methods that relies on 3DGS. Furthermore, to supplement the aforementioned methods, we have included VastGaussian (Lin et al. 2024) and MegaNeRF (Turki, Ramanan, and Satyanarayanan 2022) as additional comparative methods in the analysis of rendering outcomes.

**Datasets and Metrics** We employ GigaGS on datasets consisting of real-life aerial large-scale scenes, which encompass the *Building* and *Rubble* scenes extracted from Mill-19 (Turki, Ramanan, and Satyanarayanan 2022), along with the *Sci-Art*, *Campus*, and *Residence* scenes sourced from Urbanscene3d (Liu, Xue, and Huang 2021). To maintain consistency, we employ the same dataset partitioning as MegaNeRF (Turki, Ramanan, and Satyanarayanan 2022). We utilized visual quality metrics, namely PSNR, SSIM, and LPIPS (Zhang et al. 2018), to compare the rendering quality on the test set. Additionally, we compared the visualizations of the results on the test set with methods capable of extracting depth and normal maps.

### 5.1 Visual Quality

The figure 4 compares the rendered results of the novel perspective reconstruction method along with their corresponding depth maps and normal maps. It can be observed that our GigaGS outperforms existing surface reconstruction methods in terms of surface textures and scene geometry. The existing methods based on NeRF (Li et al. 2023; Wang et al. 2021) lack fine details and exhibit blurry and erroneous structures in image rendering. Similarly, the existing methods based on 3DGS (Chen et al. 2024a; Guédon and Lepetit 2023) are plagued by artifacts, resulting in undesirable rendering outcomes. In the table 1, we quantitatively compare the test set results of the aforementioned methods, along with the inclusion of two additional large-scale novel
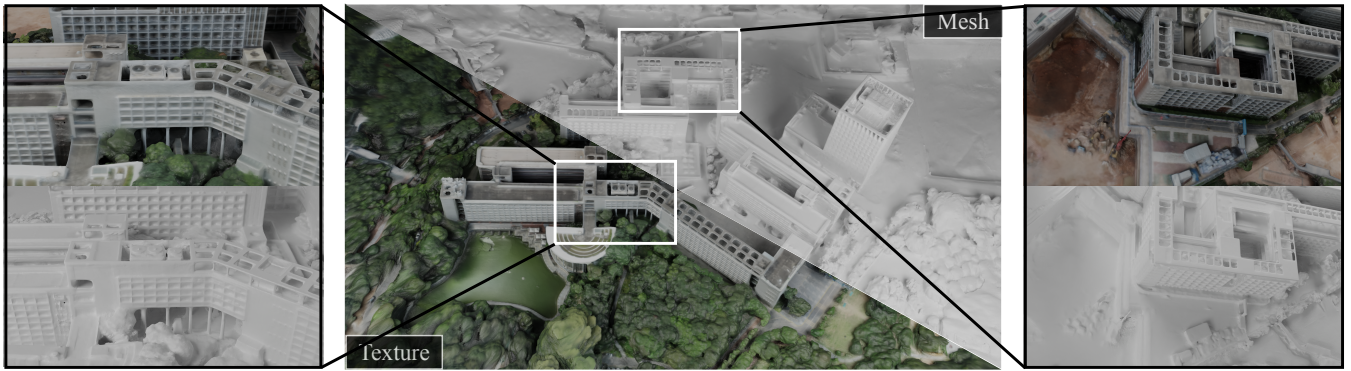
Figure 5: **Visualization of surface reconstruction.** The figures showcase the results of training GigaGS on real aerial scenes, followed by rendering multiple RGB and depth maps, and ultimately obtaining the surface reconstruction results using TSD-Fusion (Zeng et al. 2017).

Table 1: **Quantitative results of rendering quality.** We report SSIM↑, PSNR↑ and LPIPS↓ on test views. "Red", and "Yellow" denote the best and second-best results. "-" indicates that training cannot proceed due to out of memory.

|  | Building | | | Rubble | | | Campus | | | Residence | | | Sci-Art | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | SSIM | PSNR | LPIPS | SSIM | PSNR | LPIPS | SSIM | PSNR | LPIPS | SSIM | PSNR | LPIPS | SSIM | PSNR | LPIPS |
| **No mesh (except GigaGS)** | | | | | | | | | | | | | | | |
| Mega-NeRF | 0.569 | 21.48 | 0.378 | 0.575 | 24.70 | 0.407 | 0.561 | 23.93 | 0.513 | 0.648 | 22.86 | 0.330 | 0.769 | 26.25 | 0.263 |
| VastGaussian | 0.804 | 23.50 | 0.130 | 0.823 | 26.92 | 0.132 | 0.816 | 26.00 | 0.151 | 0.852 | 24.25 | 0.124 | 0.885 | 26.81 | 0.121 |
| GigaGS(ours) | 0.905 | 26.69 | 0.125 | 0.837 | 25.10 | 0.167 | 0.773 | 22.79 | 0.254 | 0.822 | 22.30 | 0.190 | 0.883 | 24.34 | 0.158 |
| **With mesh** | | | | | | | | | | | | | | | |
| PGSR | 0.480 | 16.12 | 0.573 | 0.728 | 23.09 | 0.334 | 0.399 | 14.02 | 0.721 | 0.746 | 20.57 | 0.289 | 0.799 | 19.72 | 0.275 |
| SuGaR | 0.507 | 17.76 | 0.455 | 0.577 | 20.69 | 0.453 | - | - | - | 0.603 | 18.74 | 0.406 | 0.698 | 18.60 | 0.349 |
| NeuS | 0.463 | 18.01 | 0.611 | 0.480 | 20.46 | 0.618 | 0.412 | 14.84 | 0.709 | 0.503 | 17.85 | 0.533 | 0.633 | 18.62 | 0.472 |
| Neuralangelo | 0.582 | 17.89 | 0.322 | 0.625 | 20.18 | 0.314 | 0.607 | 19.48 | 0.373 | 0.644 | 18.03 | 0.263 | 0.769 | 19.10 | 0.231 |
| GigaGS(ours) | 0.905 | 26.69 | 0.125 | 0.837 | 25.10 | 0.167 | 0.773 | 22.79 | 0.254 | 0.822 | 22.30 | 0.190 | 0.883 | 24.34 | 0.158 |

view synthesis (NVS) methods, solely for the purpose of comparing rendering quality. It can be observed that our GigaGS method significantly improves the rendering quantitative results of existing surface reconstruction methods, while achieving comparable performance to the NVS method.

## 5.2 Mesh Reconstruction

We utilized the method mentioned in the section 4.4 to extract a mesh from GigaGS. As shown in the figure 5, our approach enables the extraction of high-quality meshes while ensuring high-quality rendering. This capability holds potential to support a wide range of applications, such as navigation, simulation, and virtual reality (VR).

## 5.3 Quantity of 3D Gaussian Splatting

In the figure 6, we illustrate the quantity of 3DGS obtained through optimization in various scenarios. Owing to our stratified scene representation and partition-based optimization strategy, we are able to represent scenes with an in-

creased number of 3DGS. As the volume of scene data in practical applications grows, it can even approach the giga-level. Consequently, GigaGS can maximize the capture of scene details.

## 6 Conclusion

In this paper, we propose GigaGS. To the best of our knowledge, GigaGS is the first work for large-scale scene surface reconstruction with 3D Gaussian Splatting. Through careful design, GigaGS deliver high quality 3D surface and can process large scenes in parallel.

**Limitation** The performance of 3D Gaussian is highly correlated to the performance of COLMAP, which may degrade the performance especially for textureless regions.

## References

Charatan, D.; Li, S.; Tagliasacchi, A.; and Sitzmann, V. 2024. pixelSplat: 3D Gaussian Splats from Im-
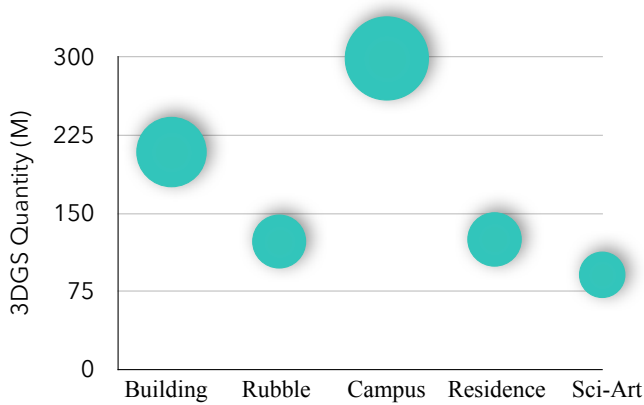
Figure 6: **Quantity of 3D Gaussian Splatting in five scenes.** This figure shows the actual number of 3DGS obtained by Our method.

age Pairs for Scalable Generalizable 3D Reconstruction. arXiv:2312.12337.

Chen, A.; Xu, Z.; Geiger, A.; Yu, J.; and Su, H. 2022. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision*, 333–350. Springer.

Chen, D.; Li, H.; Ye, W.; Wang, Y.; Xie, W.; Zhai, S.; Wang, N.; Liu, H.; Bao, H.; and Zhang, G. 2024a. PGSR: Planar-based Gaussian Splatting for Efficient and High-Fidelity Surface Reconstruction.

Chen, Z.; Wang, F.; Wang, Y.; and Liu, H. 2024b. Text-to-3D using Gaussian Splatting. arXiv:2309.16585.

Fan, Z.; Wang, K.; Wen, K.; Zhu, Z.; Xu, D.; and Wang, Z. 2023. Lightgaussian: Unbounded 3d gaussian compression with 15x reduction and 200+ fps. *arXiv preprint arXiv:2311.17245*.

Fridovich-Keil, S.; Yu, A.; Tancik, M.; Chen, Q.; Recht, B.; and Kanazawa, A. 2022. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5501–5510.

Guédon, A.; and Lepetit, V. 2023. Sugar: Surface-aligned gaussian splatting for efficient 3d mesh reconstruction and high-quality mesh rendering. *arXiv preprint arXiv:2311.12775*.

Guo, H.; Peng, S.; Lin, H.; Wang, Q.; Zhang, G.; Bao, H.; and Zhou, X. 2022. Neural 3D Scene Reconstruction with the Manhattan-world Assumption. arXiv:2205.02836.

He, X.; Chen, J.; Peng, S.; Huang, D.; Li, Y.; Huang, X.; Yuan, C.; Ouyang, W.; and He, T. 2024. Gvgen: Text-to-3d generation with volumetric representation. *arXiv preprint arXiv:2403.12957*.

Huang, B.; Yu, Z.; Chen, A.; Geiger, A.; and Gao, S. 2024a. 2D Gaussian Splatting for Geometrically Accurate Radiance Fields. *arXiv preprint arXiv:2403.17888*.

Huang, C.; Hou, Y.; Ye, W.; Huang, D.; Huang, X.; Lin, B.; Cai, D.; and Ouyang, W. 2024b. NeRF-Det++: Incorporating Semantic Cues and Perspective-aware Depth Supervision for Indoor Multi-View 3D Detection. *arXiv preprint arXiv:2402.14464*.

Jiang, Y.; Tu, J.; Liu, Y.; Gao, X.; Long, X.; Wang, W.; and Ma, Y. 2023. GaussianShader: 3D Gaussian Splatting with Shading Functions for Reflective Surfaces. *arXiv preprint arXiv:2311.17977*.

Kerbl, B.; Kopanas, G.; Leimkühler, T.; and Drettakis, G. 2023. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4): 1–14.

Li, H.; Ye, W.; Zhang, G.; Zhang, S.; and Bao, H. 2020. Saliency guided subdivision for single-view mesh reconstruction. In *2020 International Conference on 3D Vision (3DV)*, 1098–1107. IEEE.

Li, R.; Fidler, S.; Kanazawa, A.; and Williams, F. 2024. NeRF-XL: Scaling NeRFs with Multiple GPUs. arXiv:2404.16221.

Li, Z.; Müller, T.; Evans, A.; Taylor, R. H.; Unberath, M.; Liu, M.-Y.; and Lin, C.-H. 2023. Neuralangelo: High-fidelity neural surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 8456–8465.

Lin, J.; Li, Z.; Tang, X.; Liu, J.; Liu, S.; Liu, J.; Lu, Y.; Wu, X.; Xu, S.; Yan, Y.; et al. 2024. VastGaussian: Vast 3D Gaussians for Large Scene Reconstruction. *arXiv preprint arXiv:2402.17427*.

Lindenberger, P.; Sarlin, P.-E.; Larsson, V.; and Pollefeys, M. 2021. Pixel-Perfect Structure-from-Motion with Feature-metric Refinement. In *ICCV*.

Liu, X.; Ye, W.; Tian, C.; Cui, Z.; Bao, H.; and Zhang, G. 2021. Coxgraph: multi-robot collaborative, globally consistent, online dense reconstruction system. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 8722–8728. IEEE.

Liu, Y.; Xue, F.; and Huang, H. 2021. UrbanScene3D: A Large Scale Urban Scene Dataset and Simulator.

Lu, T.; Yu, M.; Xu, L.; Xiangli, Y.; Wang, L.; Lin, D.; and Dai, B. 2023. Scaffold-gs: Structured 3d gaussians for view-adaptive rendering. *arXiv preprint arXiv:2312.00109*.

Martin-Brualla, R.; Radwan, N.; Sajjadi, M. S. M.; Barron, J. T.; Dosovitskiy, A.; and Duckworth, D. 2021. NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections. arXiv:2008.02268.

Mildenhall, B.; Srinivasan, P. P.; Tancik, M.; Barron, J. T.; Ramamoorthi, R.; and Ng, R. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1): 99–106.

Ming, Y.; Ye, W.; and Calway, A. 2022. idf-slam: End-to-end rgb-d slam with neural implicit mapping and deep feature tracking. *arXiv preprint arXiv:2209.07919*.

Müller, T.; Evans, A.; Schied, C.; and Keller, A. 2022. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *ACM Trans. Graph.*, 41(4): 102:1–102:15.

Newell, M. E.; Newell, R. G.; and Sancha, T. L. 1972. A Solution to the Hidden Surface Problem. In *Proceedings of the ACM Annual Conference - Volume 1*, ACM '72, 443–450. New York, NY, USA: Association for Computing Machinery. ISBN 9781450374910.

Ren, K.; Jiang, L.; Lu, T.; Yu, M.; Xu, L.; Ni, Z.; and Dai, B. 2024. Octree-GS: Towards Consistent Real-time Rendering with LOD-Structured 3D Gaussians. *arXiv preprint arXiv:2403.17898*.

Schonberger, J. L.; and Frahm, J.-M. 2016. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 4104–4113.

Sun, C.; Sun, M.; and Chen, H.-T. 2022. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5459–5469.

Tang, J.; Chen, Z.; Chen, X.; Wang, T.; Zeng, G.; and Liu, Z. 2024a. LGM: Large Multi-View Gaussian Model for High-Resolution 3D Content Creation. arXiv:2402.05054.

Tang, J.; Ren, J.; Zhou, H.; Liu, Z.; and Zeng, G. 2023. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*.

Tang, Z.; Ye, W.; Wang, Y.; Huang, D.; Bao, H.; He, T.; and Zhang, G. 2024b. ND-SDF: Learning Normal Deflection Fields for High-Fidelity Indoor Reconstruction. *arxiv preprint*.

Turki, H.; Ramanan, D.; and Satyanarayanan, M. 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.

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

Wu, G.; Yi, T.; Fang, J.; Xie, L.; Zhang, X.; Wei, W.; Liu, W.; Tian, Q.; and Xinggang, W. 2023. 4D Gaussian Splatting for Real-Time Dynamic Scene Rendering. *arXiv preprint arXiv:2310.08528*.

Xu, Y.; Shi, Z.; Yifan, W.; Chen, H.; Yang, C.; Peng, S.; Shen, Y.; and Wetzstein, G. 2024. GRM: Large Gaussian Reconstruction Model for Efficient 3D Reconstruction and Generation. arXiv:2403.14621.

Yariv, L.; Gu, J.; Kasten, Y.; and Lipman, Y. 2021. Volume Rendering of Neural Implicit Surfaces. arXiv:2106.12052.

Ye, W.; Chen, S.; Bao, C.; Bao, H.; Pollefeys, M.; Cui, Z.; and Zhang, G. 2023a. IntrinsicNeRF: Learning Intrinsic Neural Radiance Fields for Editable Novel View Synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*.

Ye, W.; Chen, X.; Zhan, R.; Huang, D.; Huang, X.; Zhu, H.; Bao, H.; Ouyang, W.; He, T.; and Zhang, G. 2024a. DATAP-SfM: Dynamic-Aware Tracking Any Point for Robust Dense Structure from Motion in the Wild. *arxiv preprint*.

Ye, W.; Lan, X.; Chen, S.; Ming, Y.; Yu, X.; Bao, H.; Cui, Z.; and Zhang, G. 2023b. PVO: Panoptic Visual Odometry. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 9579–9589.

Ye, W.; Li, H.; Gao, Y.; Dai, Y.; Chen, J.; Dong, N.; Zhang, D.; Bao, H.; Ouyang, W.; Qiao, Y.; He, T.; and Zhang, G.

2024b. FedSurfGS: Scalable 3D Surface Gaussian Splatting with Federated Learning for Large Scene Reconstruction. *arxiv preprint*.

Ye, W.; Yu, X.; Lan, X.; Ming, Y.; Li, J.; Bao, H.; Cui, Z.; and Zhang, G. 2022. Deflowslam: Self-supervised scene motion decomposition for dynamic dense slam. *arXiv preprint arXiv:2207.08794*.

Yoo, J.-C.; and Han, T. H. 2009. Fast normalized cross-correlation. *Circuits, systems and signal processing*, 28: 819–843.

Yu, M.; Lu, T.; Xu, L.; Jiang, L.; Xiangli, Y.; and Dai, B. 2024. GSDF: 3DGS Meets SDF for Improved Rendering and Reconstruction. arXiv:2403.16964.

Yu, Z.; Chen, A.; Huang, B.; Sattler, T.; and Geiger, A. 2023. Mip-splatting: Alias-free 3d gaussian splatting. *arXiv preprint arXiv:2311.16493*.

Yu, Z.; Peng, S.; Niemeyer, M.; Sattler, T.; and Geiger, A. 2022. Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. *Advances in neural information processing systems*, 35: 25018–25032.

Yu, Z.; Sattler, T.; and Geiger, A. 2024. Gaussian Opacity Fields: Efficient and Compact Surface Reconstruction in Unbounded Scenes. arXiv:2404.10772.

Zeng, A.; Song, S.; Nießner, M.; Fisher, M.; Xiao, J.; and Funkhouser, T. 2017. 3dmatch: Learning local geometric descriptors from rgb-d reconstructions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 1802–1811.

Zhang, K.; Riegler, G.; Snavely, N.; and Koltun, V. 2020. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*.

Zhang, R.; Isola, P.; Efros, A. A.; Shechtman, E.; and Wang, O. 2018. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 586–595.