
GaussianPro: 3D Gaussian Splatting with Progressive Propagation

Kai Cheng^{*1} Xiaoxiao Long^{*2} Kaizhi Yang¹ Yao Yao³ Wei Yin⁴ Yuexin Ma⁵ Wenping Wang⁶
Xuejin Chen¹

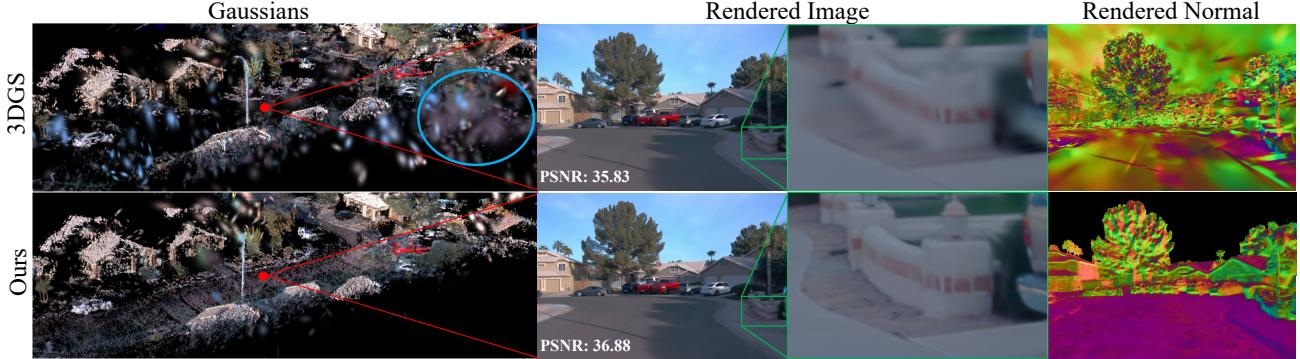


Figure 1: The sparse SfM points and less-constrained densification strategies of 3DGs pose challenges in optimizing 3D Gaussians, particularly for textureless areas. 3DGs generates incorrect Gaussians (blue circle) to be over-fitted on the training images, leading to a noticeable performance drop in novel view rendering with erroneous geometries.

Abstract

The advent of 3D Gaussian Splatting (3DGs) has recently brought about a revolution in the field of neural rendering, facilitating high-quality renderings at real-time speed. However, 3DGs heavily depends on the initialized point cloud produced by Structure-from-Motion (SfM) techniques. When tackling with large-scale scenes that unavoidably contain texture-less surfaces, the SfM techniques always fail to produce enough points in these surfaces and cannot provide good initialization for 3DGs. As a result, 3DGs suffers from difficult optimization and low-quality renderings. In this paper, inspired by classical multi-view stereo (MVS) techniques, we propose GaussianPro, a novel method that applies a progressive propagation strategy to guide the densification of the 3D Gaussians. Compared to the simple split and clone strategies used in 3DGs, our method leverages the priors of the existing reconstructed geometries of the scene and patch matching techniques to produce new Gaus-

sians with accurate positions and orientations. Experiments on both large-scale and small-scale scenes validate the effectiveness of our method, where our method significantly surpasses 3DGs on the Waymo dataset, exhibiting an improvement of 1.15dB in terms of PSNR.

1. Introduction

Novel view synthesis is an important but challenging task in computer vision and computer graphics that aims to generate images of novel viewpoints in the captured scene. It has extensive applications in various domains, including virtual reality (Deng et al., 2022b), autonomous driving (Yang et al., 2023; Cheng et al., 2023), and 3D content generation (Poole et al., 2022; Tang et al., 2023). Recently, Neural Radiance Fields (NeRF) (Mildenhall et al., 2020) has significantly boosted this task, achieving high-fidelity renderings without explicitly modeling 3D scenes, texture and illumination. However, due to the heavy manner of volume rendering, NeRFs still suffer from slow rendering speed, although various efforts (Müller et al., 2022; Barron et al., 2022; 2023; Chen et al., 2022; Xu et al., 2022) have been made.

To achieve real-time neural rendering, 3D Gaussian Splatting (3DGs) (Kerbl et al., 2023) has been developed. It models the scenes explicitly as 3D Gaussians with learnable attributes and performs rasterization of the Gaus-

^{*}Equal contribution ¹University of Science and Technology of China ²The University of Hong Kong ³Nanjing University ⁴The University of Adelaide ⁵ShanghaiTech University
⁶Texas A&M University. Correspondence to: Xuejin Chen <xjchen99@ustc.edu.cn>.

sians to produce renderings. The splatting strategy avoids time-consuming ray sampling and allows parallel computations, thus yielding high efficiency and fast rendering. 3DGS heavily relies on the sparse point clouds produced by Structure-from-Motion (SfM) techniques to initialize the Gaussians, e.g., their positions, colors, and shapes. Moreover, clone and split strategies to create more new Gaussians to achieve a complete coverage of the scene. However, the densification strategies to create 3D Gaussians lead to two main limitations. 1) **Sensitive to Gaussian Initialization.** The SfM techniques always fail to produce 3D points and leave empty in textureless regions, and therefore the densification strategy struggles to generate reliable Gaussians to cover the scene with a poor initialization. 2) **Ignore the priors of the existing reconstructed geometry.** The new Gaussians are either cloned to be the same as the old Gaussians or initialized with random positions and orientations. The less-constrained densification leads to difficulties in the optimization of 3D Gaussians, e.g., noisy geometries, and few Gaussians in texture-less regions, finally degrading the rendering quality. As shown in Figure 1, the results of 3DGS contain many noisy Gaussians and some regions are not covered by enough Gaussians.

In this paper, we propose a novel progressive propagation strategy to facilitate 3DGS, which could produce more compact and accurate 3D Gaussians and therefore improve the rendering quality, especially in texture-less surfaces. The key idea of our method is to fully leverage the reconstructed scene geometries as priors and classical patch matching techniques to progressively produce new Gaussians with accurate positions and orientations.

Specifically, we consider Gaussian densification in both 3D world space and 2D image space. For each input image, we render the depth and normal map by accumulating the positions and orientations of 3D Gaussians via alpha blending. Based on the observation that the neighboring pixels are likely to share similar depth and normal values, for a pixel, we iteratively propagate the depth and normal values of its neighboring pixels to this pixel to formulate a set of candidates. With the candidates, we leverage classical patch matching techniques to pick up the best candidate that satisfies the multi-view photometric consistency constraint, thus yielding new depth and normal for each pixel (named as propagated depth/normal). We select the pixels whose propagated depth is significantly different from the rendered depth since large differences imply that the existing 3D Gaussians may not accurately capture the true geometry. As a result, we explicitly back-project the selected pixels using the propagated depths into 3D space and initialize them as new Gaussians. Additionally, we leverage the propagated normals to regularize the orientations of 3D Gaussians, further improving the reconstructed 3D geometry and rendering quality.

Our proposed progressive propagation strategy could produce more compact and accurate 3D Gaussians by transferring of accurate geometric information from well-modeled regions to under-modeled regions. As shown in Figure 1, compared to 3DGS, our method produces more accurate and compact Gaussians and therefore achieves a better coverage of the 3D scene. Experiments on public datasets such as Waymo and MipNeRF360 validate that our proposed strategy significantly boosts the performance of 3DGS. Overall, the contributions of our method are summarized as:

- We propose a novel Gaussian propagation strategy that guides the densification to produce more compact and accurate Gaussians, particularly in low-texture regions.
- We additionally leverage a planar loss that provides a further constraint in the optimization of Gaussians.
- Our method achieves new state-of-the-art rendering performance on the Waymo and MipNeRF360 datasets. Our method also presents robustness to the varying numbers of input images.

2. Related Work

2.1. Multi-view Stereo

Multi-view stereo (MVS) aims to reconstruct a 3D model from a collection of posed images, which can be further combined with traditional rendering algorithms to generate novel views. Traditional methods (Campbell et al., 2008; Furukawa & Ponce, 2009; Bleyer et al., 2011; Furukawa et al., 2015; Schönberger et al., 2016; Xu & Tao, 2019) explicitly establish pixel correspondences between images based on hand-crafted image features and then optimize the 3D structure to achieve the best pixel correspondences among images. Learning-based MVS methods (Yao et al., 2018; Vakalopoulou et al., 2018; Long et al., 2020; Chen et al., 2019; Long et al., 2021; Ma et al., 2022; Long et al., 2022; Feng et al., 2023) implicitly build multi-view correspondences with learnable features and regress depth or 3D volume based on the features in an end-to-end framework. In this paper, we draw inspiration from depth optimization in MVS to improve the geometry of the Gaussians, thereby achieving better rendering results.

2.2. Neural Radiance Field

NeRF combines deep learning techniques with the 3D volumetric representation, transforming a 3D scene into a learnable continuous density field. Utilizing ray marching in volume rendering, NeRF is able to achieve high-quality novel view synthesis without explicit modeling of the 3D scene and illumination. To further improve the

rendering quality, some approaches (Barron et al., 2021; Xu et al., 2022; Barron et al., 2023) directly improve the point sampling strategy in ray marching for more accurate modeling of the volume rendering process. Others (Barron et al., 2022; Wang et al., 2023) improve rendering by reparameterizing the scene to generate more compact scene representation and easier learning process. Additionally, regularization terms (Deng et al., 2022a; Yu et al., 2022) could be introduced to constrain the scene representation towards a closer approximation of real geometry. Despite these advancements, NeRF still incurs high computational costs during rendering. Since NeRF employs MLPs to represent the scene, the computation and optimization of any point in the scene are dependent on the entire MLP. Many works propose novel scene representations to accelerate rendering. They replace MLPs with sparse voxels (Liu et al., 2020; Fridovich-Keil et al., 2022), hash tables (Müller et al., 2022), or triplane (Chen et al., 2022), allowing the computation and optimization of each point to be localized to the corresponding local region of the scene. Although these methods significantly improve rendering speed, real-time rendering is still challenging due to the inherent ray marching strategy in volume rendering.

2.3. 3D Gaussian Splatting

3DGS employs a spatting-based rasterization (Zwicker et al., 2002) approach to project anisotropic 3D Gaussians onto a 2D screen. It computes the pixel’s color by performing depth sorting and α -blending on the projected 2D Gaussians, which avoids the sophisticated sampling strategy of ray marching and achieves real-time rendering. Some concurrent works have made improvements to 3DGS. Firstly, 3DGS is sensitive to sampling frequency, i.e., changing the camera’s focal length or camera distance could result in rendering artifacts. These artifacts are addressed by introducing low-pass filtering (Yu et al., 2023) or multi-scale Gaussian representations (Yan et al., 2023). Additionally, 3DGS excessively grows Gaussians without explicitly constraining the scene’s real geometric structure, resulting in numerous redundant Gaussians and significant memory consumption. Some methods evaluate the contribution of Gaussians to rendering by their scales (Lee et al., 2023) or calculating their visibility in views (Fan et al., 2023), forcing the removal of Gaussians with small contributions. Others compress the storage of Gaussian attributes by quantization technique (Navaneet et al., 2023) or interpolating Gaussian attributes from structured grid features (Morgenstern et al., 2023; Lu et al., 2023).

Although these methods significantly reduce the storage overhead of Gaussians, they do not explicitly constrain the geometry of the Gaussians. 3DGS could grow in locations far from the real surfaces to fit different training views, resulting in redundancy and a decrease in rendering quality

for new viewpoints. This paper considers the planar prior in the scene, explicitly constraining the growth of Gaussians close to the real surfaces. This approach enables Gaussians to better fit the real geometry of the scene, achieving improved rendering and more compact representation.

3. Preliminaries

3DGS (Kerbl et al., 2023) models the 3D scene as a set of anisotropic 3D Guassians, which are further rendered to images using the splatting-based rasterization technique (Zwicker et al., 2002). For each 3D Gaussian G , it is defined as:

$$G(\mathbf{x}) = e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu})}, \quad (1)$$

where $\boldsymbol{\mu} \in \mathbb{R}^{3 \times 1}$ refers to its mean vector, $\boldsymbol{\Sigma} \in \mathbb{R}^{3 \times 3}$ refers to its covariance matrix. In order to ensure the positive semi-definite property of the covariance matrix during the optimization, it is further expressed as $\boldsymbol{\Sigma} = \mathbf{R} \mathbf{S} \mathbf{S}^T \mathbf{R}^T$, where the rotation matrix $\mathbf{R} \in \mathbb{R}^{3 \times 3}$ is orthogonal, and the scale matrix $\mathbf{S} \in \mathbb{R}^{3 \times 3}$ is diagonal.

To render an image from a given viewpoint, the color of each pixel \mathbf{p} is calculated by blending N ordered Gaussians $\{G_i \mid i = 1, \dots, N\}$ overlapping \mathbf{p} as

$$\mathbf{c}(\mathbf{p}) = \sum_{i=1}^N \mathbf{c}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (2)$$

where α_i is obtained by evaluating a projected 2D Gaussian (Zwicker et al., 2002) from G_i in \mathbf{p} multiplied with a learned opacity of G_i , and \mathbf{c}_i is the learnable color of G_i . Gaussians that cover \mathbf{p} are sorted in ascending order of their depths under the current viewpoint. Through differentiable rendering techniques, all attributes of the Gaussians could be optimized end-to-end via training view reconstruction.

In order to accurately represent the scene geometry, 3DGS also employs a densification strategy to generate new Gaussians. For each training iteration, if the gradient backpropagated from the rendering loss to the current Gaussian exceeds a certain threshold, 3DGS considers that it does not sufficiently represent the corresponding 3D region. If the covariance of the Gaussian is large, it is split into two Gaussians. Conversely, if the covariance is small, it is cloned. This strategy encourages 3DGS to increase the number of Gaussians to cover the captured scene.

4. Method

4.1. Overview

In this paper, we propose a novel progressive propagation strategy to explicitly generate 3D Gaussians with accu-

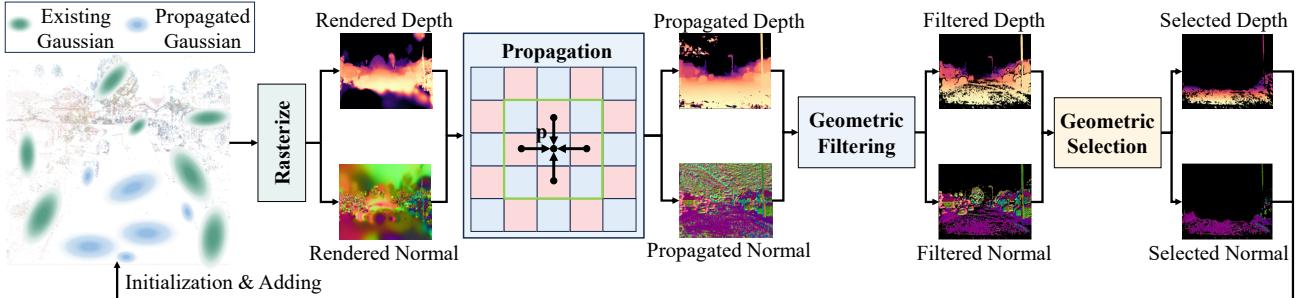


Figure 2: Progressive Propagation of Gaussian. Firstly, we render the depth and normal maps from the 3D Gaussians. Then we iteratively perform propagation operations on the rendered depths and normals to generate new depth and normal values (denoted as propagated depth and propagated normal) via patch matching techniques. We filter out the unreliable propagated depths and normals using geometric consistency, yielding filtered depths and filtered normals. Finally, we identify the regions where their rendered depths and normals significantly deviate from the filtered ones, indicating that existing Gaussians may not accurately capture the geometry and therefore need more Gaussians. Pixels in these regions are projected into the 3D space to initialize new Gaussians using the filtered depth and normal.

rate positions and orientations, thereby improving rendering quality and compactness. First, instead of only coupling with 3D space, we propose to tackle this problem in both 3D space and 2D image space. We project 3D Gaussians onto the 2D space to generate depth and normal maps, which are used to guide the growth of Gaussians (Sec. 4.2). Then we iteratively update each pixel’s depth and normal based on the propagated ones from its neighboring pixels. Pixels whose new depth is significantly different from the initial depth are projected back to 3D space as 3D points and these points are further initialized as new Gaussians (Sec. 4.3). Additionally, a planar loss is also incorporated to further regularize the geometry of the Gaussians, yielding more accurate geometries (Sec. 4.4). The overall training strategy is introduced in Sec. 4.5.

4.2. Hybrid Geometric Representation

In this section, we propose a hybrid Geometric representation that combines 3D Gaussians with 2D view-dependent depth and normal maps, where the 2D representations are utilized to assist the densification of Gaussians.

Due to the discrete and irregular topology of the 3D Gaussians, it is inconvenient to perceive the connectivity of geometries, like searching neighboring Gaussians on a local surface. As a result, it’s difficult to perceive the existing geometry to guide the Gaussian densification. Inspired by the classical MVS methods, we propose to tackle this challenge by mapping the 3D Gaussians into structured 2D image space. This mapping allows us to efficiently determine the neighbors of the Gaussians and propagate geometric information among them. Specifically, when Gaussians are located on the same local plane in 3D space, their 2D projections should also be in adjacent regions and exhibit sim-

ilar geometric properties, i.e. depth and normal.

The depth value of Gaussian. For each viewpoint with camera extrinsics $[\mathbf{W}, \mathbf{t}] \in \mathbb{R}^{3 \times 4}$, the center μ_i of a Gaussian G_i can be projected into the camera coordinate system as μ'_i :

$$\mu'_i = \begin{bmatrix} x_i \\ y_i \\ z_i \end{bmatrix} = \mathbf{W}\mu_i + \mathbf{t}, \quad (3)$$

where z_i refers to the Gaussian’s depth under the current viewpoint.

The normal value of Gaussian. In a Gaussian G_i , its covariance matrix is formulated as $\Sigma_i = \mathbf{R}_i \mathbf{S}_i \mathbf{S}_i^T \mathbf{R}_i^T$. The rotation matrix \mathbf{R}_i determines its three orthogonal eigenvectors while scaling matrix $\mathbf{S}_i \in \mathbb{R}^{3 \times 3}$ determines the scale along the eigenvector directions. The covariance matrix Σ_i of a 3D Gaussian could be compared to a representation of the shape of an ellipsoid, where the eigenvectors correspond to ellipsoid’s axes and the scales refer to the lengths of the axes. According to the Gaussian-Shader (Jiang et al., 2023), the Gaussian sphere gradually becomes flattened and approaches a plane during the optimization process. Therefore, the direction of its shortest axis can approximate the normal direction \mathbf{n}_i of the Gaussian, which is induced by

$$\mathbf{n}_i = \mathbf{R}_i[r, :], r = \text{argmin} ([s_1, s_2, s_3]), \quad (4)$$

where $\text{diag}(s_1, s_2, s_3) = \mathbf{S}_i$, $\text{argmin}(\cdot)$ is the operation to find the index of the minimum value.

Finally, the 2D depth and normal map under the current viewpoint are rendered based on α -blending defined in Eq. 2, where the attribute color \mathbf{c}_i is replaced by Gaussian’s depth z_i and normal \mathbf{n}_i .

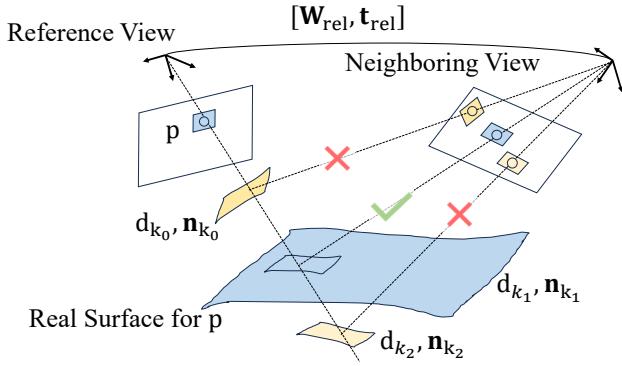


Figure 3: Patch matching. To select the best plane candidate for pixel p during propagation, we perform homography transformation between p and each plane candidate, thus yielding the possible corresponding pixels of the neighboring view. The plane candidate that exhibits the highest color consistency between p and its possible paired pixel is chosen to be the solution. The chosen plane candidate is used to update the depth and normal of pixel p .

4.3. Progressive Gaussian Propagation

In this section, we introduce the progressive Gaussian propagation strategy, which allows for propagating accurate geometry from well-modeled regions to under-modeled regions, enabling producing new Gaussians. As shown in Figure 2, with the rendered depth maps and normal maps, we employ patch matching (Barnes et al., 2009) to propagate the depth and normal information from neighboring pixels to the current pixel, which produces new depths and normals (named as propagated depth/normal). We further perform geometric filtering and selection operations to pick up the pixels that need more Gaussians and leverage their propagated depths and normals to initialize new Gaussians.

Plane Definition. To achieve the propagation, the depth and normal of each pixel need to be converted to a 3D local plane first. For each pixel with its coordinate \mathbf{p} , the corresponding 3D local plane is parameterized as (d, \mathbf{n}) , where \mathbf{n} is the pixel’s rendered normal, and d is the distance from the origin of the camera coordinate to the local plane calculated as:

$$d = z\mathbf{n}^\top \mathbf{K}^{-1}\tilde{\mathbf{p}}, \quad (5)$$

where $\tilde{\mathbf{p}}$ is the homogeneous coordinate of \mathbf{p} , z is pixel’s rendered depth, and \mathbf{K} refers to the camera intrinsic.

Candidate Selection. After defining the 3D local plane, the neighbors of each pixel need to be selected for propagation. We follow the checkerboard pattern defined in ACMH (Xu & Tao, 2019) to select neighboring pixels. For clarity, we illustrate the propagation of a pixel with its four nearest pixels. For each pixel, a set of plane candidates $\{(d_{k_l}, \mathbf{n}_{k_l}) \mid l \in \{0, 1, 2, 3, 4\}\}$ is obtained through propagation (k_l refers to the index of pixel p and its four neighbors).

boring pixels).

Patch Matching. After obtaining the plane candidates, the optimal plane for each pixel is determined through patch matching. For a pixel p with its coordinate \mathbf{p} , a homography transformation \mathbf{H} is performed based on each plane candidate $(d_{k_l}, \mathbf{n}_{k_l})$, which warps \mathbf{p} to \mathbf{p}' in the neighboring frame as:

$$\tilde{\mathbf{p}}' \simeq \mathbf{H}\tilde{\mathbf{p}}, \quad (6)$$

where $\tilde{\mathbf{p}}'$ is the homogeneous coordinate of \mathbf{p}' , and \mathbf{H} can be induced as:

$$\mathbf{H} = \mathbf{K} \left(\mathbf{W}_{\text{rel}} - \frac{\mathbf{t}_{\text{rel}} \mathbf{n}_{k_l}^\top}{d_{k_l}} \right) \mathbf{K}^{-1}, \quad (7)$$

where $[\mathbf{W}_{\text{rel}}, \mathbf{t}_{\text{rel}}]$ is the relative transformation from the reference view to the neighboring view. Finally, the color consistency of p and p' is evaluated based on NCC (Normalized Cross Correlation) (Yoo & Han, 2009). The local plane of p will be updated to the plane candidate with the best color consistency. Fig. 3 also provides an intuitive visualization of this process. The propagation for plane candidates is iterated u times to transmit effective geometric information over a large region. Then the pixel’s depth and normal are updated from the propagated plane, ultimately resulting in the propagated depth and normal maps in Fig. 2.

Geometric Filtering and Selection. Due to the inevitable errors in the propagated results, we filter out inaccurate depth and normal through multi-view geometric consistency check (Schönberger et al., 2016) and obtain filtered depth and normal maps. Finally, we calculate the absolute relative difference between the filtered depth and rendered depth. For regions with an absolute relative difference greater than the threshold σ , we consider that existing Gaussians fail to model these regions accurately. Therefore, we project pixels in these regions back to the 3D space and initialize them as 3D Gaussians using the same initialization in 3DGS. These Gaussians are then added to the existing Gaussians for further optimization.

4.4. Plane Constraint Optimization

In the original 3DGS, the optimization only relies on image reconstruction loss without incorporating any geometric constraints. As a result, the optimized Gaussian shapes may deviate significantly from the actual surface geometry. This deviation leads to a decline in the rendering quality when viewed from a new viewpoint, particularly for large-scale scenes with limited views. As shown in Fig. 4, the shape of Gaussians in 3DGS differs significantly from the road’s geometry, resulting in severe rendering artifacts when viewed from a novel viewpoint. In this section, we propose a planar constraint that encourages the shape of Gaussians to closely resemble the real surface.

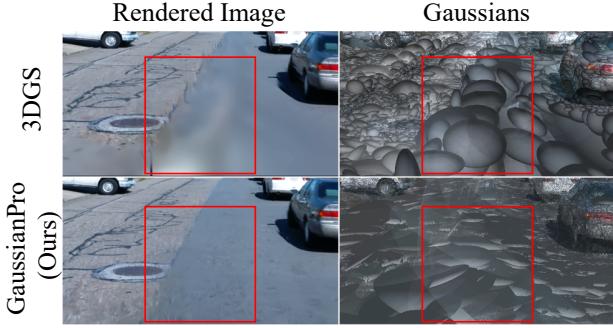


Figure 4: Visual comparisons with 3DGS on novel view synthesis. The rendered image of 3DGS contains severe artifacts since the Gaussian spheres are out of order and do not accurately model the true geometry. On the contrary, our method faithfully captures the details of the road, and its Gaussian spheres are more compact and orderly.

Specifically, the propagated 2D normal map in Section 4.3 represents the orientation of the planes in the scene. We explicitly enforce the consistency between Gaussian’s rendered normal and the propagated normal with \mathcal{L}_1 and angular loss as $\mathcal{L}_{\text{normal}}$:

$$\mathcal{L}_{\text{normal}} = \sum_{\mathbf{p} \in \mathcal{Q}} \|\hat{\mathbf{N}}(\mathbf{p}) - \bar{\mathbf{N}}(\mathbf{p})\|_1 + \left\| 1 - \hat{\mathbf{N}}(\mathbf{p})^\top \bar{\mathbf{N}}(\mathbf{p}) \right\|_1, \quad (8)$$

where $\hat{\mathbf{N}}$ is the rendered normal map, $\bar{\mathbf{N}}$ is the propagated normal map, and \mathcal{Q} refers to the set of valid pixels after the geometric filtering in Section 4.3.

Additionally, to ensure that the shortest axis of the Gaussian could represent the normal direction, we incorporate a scale regularization loss $\mathcal{L}_{\text{scale}}$ in NeuSG (Chen et al., 2023). This loss constrains the minimum scale in Gaussian to be close to zero, effectively flattening the Gaussians towards a planar shape. Finally, the plane constraint can be expressed as the weighted sum of two losses:

$$\mathcal{L}_{\text{planar}} = \beta \mathcal{L}_{\text{normal}} + \gamma \mathcal{L}_{\text{scale}}. \quad (9)$$

4.5. Training Strategy

In summary, we incorporate the progressive Gaussian propagation strategy into 3DGS, activating it every m iterations in the optimization, where we set $m = 50$. The propagated normal maps are saved for computing the planar constraint loss. Our final training loss \mathcal{L} consists of the image reconstruction loss \mathcal{L}_1 and $\mathcal{L}_{\text{D-SSIM}}$ in 3DGS with the proposed planar constraint loss, as illustrated in Eq. 10.

$$\mathcal{L} = (1 - \lambda) \mathcal{L}_1 + \lambda \mathcal{L}_{\text{D-SSIM}} + \mathcal{L}_{\text{planar}}, \quad (10)$$

where the weight λ is set to 0.2 same with 3DGS.

5. Experiment

5.1. Datasets and Implementation Details

Datasets. We conduct our experiments in a large-scale urban dataset Waymo (Sun et al., 2020), and the common NeRF benchmark Mip-NeRF360 dataset. (Caesar et al., 2020). On the Waymo dataset, we randomly select nine scenes for evaluation. To evaluate the performance of novel view synthesis, following the common settings, we select one of every eight images as testing images and the remaining ones as training data. We apply the three widely-used metrics for evaluation, *i.e.*, peak signal-to-noise ratio (PSNR), structural similarity index measure (SSIM), and the learned perceptual image patch similarity (LPIPS) (Zhang et al., 2018).

Implementation Details. Our method is built upon the popular open-source 3DGS code base (Kerbl et al., 2023). In alignment with the approach described in 3DGS, our models are trained for 30,000 iterations across all scenes following 3DGS’s training schedule and hyperparameters. Besides the original clone and split Gaussian densification strategies used in 3DGS, We additionally perform our proposed progressive propagation strategy every 50 training iterations where propagation is performed 3 times. The threshold σ of the absolute relative difference is set to 0.8. For the planar loss, we set $\beta = 0.001$ and $\gamma = 100$. All experiments are conducted on an RTX 3090 GPU.

5.2. Quantitative and Qualitative Results

As shown in Tab. 1, we compare our method with the state-of-the-art (SOTA) methods, including Instant-NGP (Müller et al., 2022), Mip-NeRF360 (Barron et al., 2022), ZipNeRF (Barron et al., 2023), and 3DGS (Kerbl et al., 2023).

Results on Waymo. On the large-scale urban dataset Waymo, our method significantly outperforms others in all evaluation metrics. Due to the presence of textureless regions in street views, initializing point clouds in these regions becomes a challenge for SfM. Consequently, it is difficult for 3DGS to densify Gaussians that accurately represent the geometry of the scene in these regions. Otherwise, our propagation strategy accurately complements the missing geometry in the scene. Additionally, our planar constraint allows for better modeling of the scene’s planes. Therefore, compared to the baseline 3DGS, our method significantly improves PSNR by 1.15 dB. The visual results presented in Fig. 5 show that our method achieves sharp details and better renderings in both rich-texture and texture-less regions.

Results on MipNeRF360. On the MipNeRF360, we retrain 3DGS using our generated SfM point clouds since we observed the SfM points used in their official code base can be improved. We report the quantitative results of

Table 1: Quantitative comparisons on Waymo and MipNeRF-360 datasets. We indicate the best and second best with bold and underlined respectively. 3DGS* refers to the results obtained by 3DGS retrained with better SfM point clouds.

| Method | FPS \uparrow | PSNR \uparrow | Waymo | | MipNeRF 360 | | |
|------------------------------------|----------------|-----------------|-----------------|--------------------|-----------------|-----------------|--------------------|
| | | | SSIM \uparrow | LPIPS \downarrow | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow |
| Instant-NGP (Müller et al., 2022) | 3 | 30.98 | 0.886 | 0.281 | 25.59 | 0.699 | 0.331 |
| Mip-NeRF 360 (Barron et al., 2022) | 0.02 | 30.09 | 0.909 | 0.262 | 27.69 | 0.792 | 0.237 |
| Zip-NeRF (Barron et al., 2023) | 0.09 | <u>34.22</u> | <u>0.939</u> | <u>0.205</u> | 28.54 | 0.828 | 0.189 |
| 3DGS* (Kerbl et al., 2023) | <u>103</u> | 33.53 | 0.938 | 0.226 | 27.21 | 0.815 | 0.214 |
| 3DGS (Retrained) | 102 | - | - | - | 27.88 | 0.824 | 0.209 |
| GaussianPro (Ours) | 108 | 34.68 | 0.949 | 0.191 | 27.92 | 0.825 | 0.208 |

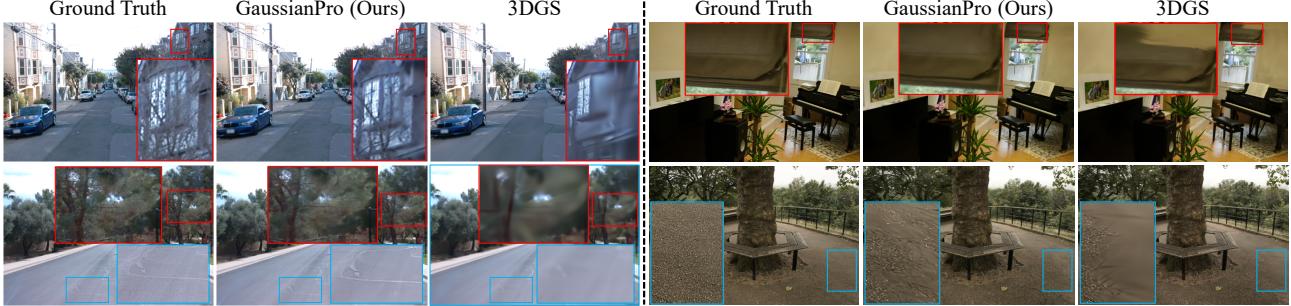


Figure 5: Rendering results on the Waymo (left) and MipNeRF360 (right) datasets. Compared to 3DGS, we have achieved a noticeable improvement in both texture-less surfaces and sharp details.

Table 2: Ablation study on the proposed propagation strategy and planar constraint.

| Propagation | Planar | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow |
|-------------|----------|-----------------|-----------------|--------------------|
| \times | \times | 33.53 | 0.938 | 0.226 |
| \times | ✓ | 34.02 | 0.942 | 0.218 |
| ✓ | \times | 34.48 | 0.946 | 0.203 |
| ✓ | ✓ | 34.68 | 0.949 | 0.191 |

the original 3DGS (denoted as 3DGS*) and our retrained 3DGS in Tab. 1. Our method achieves comparable results with 3DGS with a slight improvement. The MipNeRF360 dataset contains quite small-scale natural and indoor scenes with rich textures, so the SfM techniques usually provide a high-quality point cloud for initialization and the simple clone and split densification strategies don’t show a bottleneck in the small-scale scenes. For indoor scenes with some weak-texture surfaces, our method still shows improvement. We report results for each scene under MipNeRF360 in the appendix to further support our conclusions. As shown in Fig 5, compared to 3DGS, our method achieves more accurate renderings and clear details.

5.3. Ablation Study

Effectiveness of the Propagation Strategy and Planar Constraint. We validate the effectiveness of the proposed propagation strategy and planar constraint in the Waymo dataset. As shown in Tab. 2, the progressive propaga-



Figure 6: Visualization of Gaussians in the *Room* scene of MipNeRF360 dataset. Our method contains fewer noisy Gaussians and achieves a more compact representation.

tion strategy (the third row) brings significant improvement compared with the baseline. This improvement can be attributed to its ability to refine the geometric representation of the scene, particularly in regions where the initial 3DGS exhibits significant errors (shown in the first and second rows of Fig. 7). The planar constraint can further enhance the rendering quality by accurately modeling the normals of the planes, as shown in the third row of Fig. 7.

The Robustness against Sparse Training Images. As the number of training images decreases, the rendering quality of neural rendering methods, including 3DGS, tends to decline. In Table 3, we present the results of training 3DGS and our method using randomly selected subsets comprising 30%, 50%, 70%, 100% of the training images from a scene in MipNeRF360. Remarkably, our method consistently achieves superior rendering results compared to 3DGS across different percentages of training images.

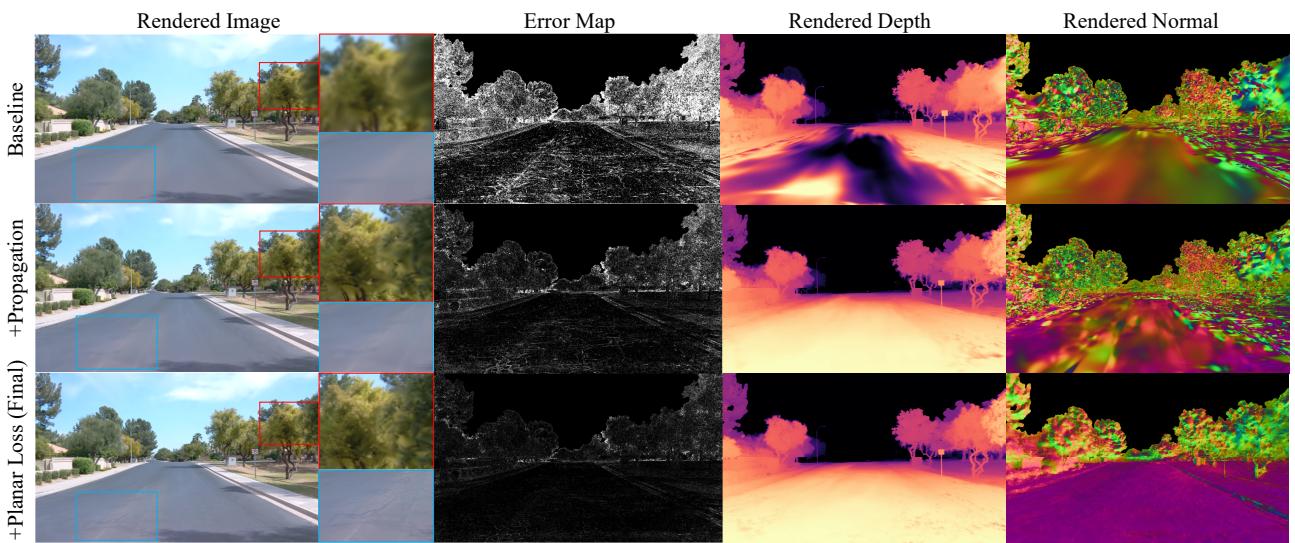


Figure 7: The progressive propagation strategy effectively enhances the geometry of the scene, resulting in improved rendering quality. The planar constraint further improves the geometry and rendering of planes.

Table 3: Comparison of 3DGS and ours with different training view ratios in the *room* scene of the MipNeRF360 dataset.

| Method | 30% | | | 50% | | | 70% | | | 100% | | |
|-------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | PSNR | SSIM | LPIPS |
| 3DGS | 28.45 | 0.896 | 0.216 | 29.97 | 0.912 | 0.203 | 30.87 | 0.921 | 0.194 | 31.71 | 0.919 | 0.192 |
| GaussianPro(Ours) | 28.64 | 0.900 | 0.210 | 30.27 | 0.914 | 0.199 | 30.93 | 0.924 | 0.189 | 31.98 | 0.927 | 0.192 |

Table 4: Efficiency analysis. We analyze the effects of initialization using SfM points or MVS points.

| Scene | Strategy | PSNR | Gaussians | Training | FPS |
|--------|------------------------|--------------|--------------|--------------|------------|
| Street | SfM points+3DGS | 35.05 | 665k | 40min | 119 |
| | MVS points+3DGS | 36.13 | 1705k | 250min | 75 |
| | SfM points+GaussianPro | 36.08 | 991k | 56min | 108 |
| Room | SfM points+3DGS | 31.71 | 1537k | 59min | 105 |
| | MVS points+3DGS | 32.05 | 1832k | 270min | 90 |
| | SfM points+GaussianPro | 31.98 | 1461k | 70min | 113 |

Efficiency Analysis. We select two typical outdoor and indoor scenes to compare the efficiency of our method with 3DGS, as shown in Tab. 4. We achieve a noticeable improvement in rendering quality with a slight increase in training time. In the case of the street scene, 3DGS uses large incorrect Gaussians to represent the ground, as shown in the blue circle of Fig. 1, resulting in fewer Gaussians compared to our method. However, for the room scene, our method results in more compact Gaussians with less noise (also shown in Fig. 6). Additionally, our method achieves a comparable real-time rendering frame rate as 3DGS.

Comparison to MVS Inputs. As our method achieves better rendering quality by improving Gaussians’ geometry, it raises the question of whether a similar effect can be achieved by directly inputting denser and more accurate MVS point clouds into 3DGS. To investigate this, we compare the results of optimizing 3DGS with the dense point cloud generated by the MVS method (Schönberger et al.,

2016). Tab. 4 shows that directly inputting the MVS point cloud significantly increases the training time (approximately 4 times) due to the additional MVS process and the large number of initial Gaussians. Moreover, the number of Gaussians increases significantly, and the rendering speed noticeably decreases, despite a slight improvement in rendering quality. Contrarily, our method achieves a favorable balance between rendering quality and efficiency.

6. Conclusion

In this paper, we propose GaussianPro, a novel progressive propagation strategy to guide Gaussian densification according to the surface structure of the scene. Based on the propagation process, we additionally introduce the plane constraints during optimization to encourage the Gaussians to better model planar surfaces. Our method demonstrates superior rendering results compared to 3DGS on both Waymo and MipNeRF360 datasets, while maintaining compact Gaussian representations. Our method shows significant improvements in structured scenes and remains robust to variations in the number of training images. However, our method does not specially model dynamic objects, and will present artifacts on these regions like all the static Gaussian methods. In the future, the recent dynamic Gaussian techniques can be incorporated into our method as complementary components to handle dynamic objects.

References

- Barnes, C., Shechtman, E., Finkelstein, A., and Goldman, D. B. Patchmatch: A randomized correspondence algorithm for structural image editing. *ACM Trans. Graph.*, 28(3):24, 2009.
- Barron, J. T., Mildenhall, B., Tancik, M., Hedman, P., Martin-Brualla, R., and Srinivasan, P. P. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 5855–5864, 2021.
- Barron, J. T., Mildenhall, B., Verbin, D., Srinivasan, P. P., and Hedman, P. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5470–5479, 2022.
- Barron, J. T., Mildenhall, B., Verbin, D., Srinivasan, P. P., and Hedman, P. Zip-nerf: Anti-aliased grid-based neural radiance fields. *Proceedings of the IEEE International Conference on Computer Vision*, 2023.
- Bleyer, M., Rhemann, C., and Rother, C. Patchmatch stereo-stereo matching with slanted support windows. In *Bmvc*, volume 11, pp. 1–11, 2011.
- Caesar, H., Bankiti, V., Lang, A. H., Vora, S., Lioung, V. E., Xu, Q., Krishnan, A., Pan, Y., Baldan, G., and Beijbom, O. nuscenes: A multimodal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11621–11631, 2020.
- Campbell, N. D., Vogiatzis, G., Hernández, C., and Cipolla, R. Using multiple hypotheses to improve depth-maps for multi-view stereo. In *Computer Vision–ECCV 2008: 10th European Conference on Computer Vision, Marseille, France, October 12–18, 2008, Proceedings, Part I 10*, pp. 766–779. Springer, 2008.
- Chen, A., Xu, Z., Geiger, A., Yu, J., and Su, H. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision*, pp. 333–350. Springer, 2022.
- Chen, H., Li, C., and Lee, G. H. Neusg: Neural implicit surface reconstruction with 3d gaussian splatting guidance. *arXiv preprint arXiv:2312.00846*, 2023.
- Chen, R., Han, S., Xu, J., and Su, H. Point-based multi-view stereo network. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 1538–1547, 2019.
- Cheng, K., Long, X., Yin, W., Wang, J., Wu, Z., Ma, Y., Wang, K., Chen, X., and Chen, X. Uc-nerf: Neural radiance field for under-calibrated multi-view cameras in autonomous driving. *arXiv preprint arXiv:2311.16945*, 2023.
- Deng, K., Liu, A., Zhu, J.-Y., and 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, 2022a.
- Deng, N., He, Z., Ye, J., Duinkharjav, B., Chakravarthula, P., Yang, X., and Sun, Q. Fov-nerf: Foveated neural radiance fields for virtual reality. *IEEE Transactions on Visualization and Computer Graphics*, 28(11):3854–3864, 2022b.
- Fan, Z., Wang, K., Wen, K., Zhu, Z., Xu, D., and Wang, Z. Lightgaussian: Unbounded 3d gaussian compression with 15x reduction and 200+ fps. *arXiv preprint arXiv:2311.17245*, 2023.
- Feng, Z., Yang, L., Guo, P., and Li, B. Cvrecon: Rethinking 3d geometric feature learning for neural reconstruction. *arXiv preprint arXiv:2304.14633*, 2023.
- Fridovich-Keil, S., Yu, A., Tancik, M., Chen, Q., Recht, B., and Kanazawa, A. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5501–5510, 2022.
- Furukawa, Y. and Ponce, J. Accurate, dense, and robust multiview stereopsis. *IEEE transactions on pattern analysis and machine intelligence*, 32(8):1362–1376, 2009.
- Furukawa, Y., Hernández, C., et al. Multi-view stereo: A tutorial. *Foundations and Trends® in Computer Graphics and Vision*, 9(1-2):1–148, 2015.
- Jiang, Y., Tu, J., Liu, Y., Gao, X., Long, X., Wang, W., and Ma, Y. Gaussianshader: 3d gaussian splatting with shading functions for reflective surfaces. *arXiv preprint arXiv:2311.17977*, 2023.
- Kerbl, B., Kopanas, G., Leimkühler, T., and Drettakis, G. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), 2023.
- Lee, J. C., Rho, D., Sun, X., Ko, J. H., and Park, E. Compact 3d gaussian representation for radiance field. *arXiv preprint arXiv:2311.13681*, 2023.
- Liu, L., Gu, J., Zaw Lin, K., Chua, T.-S., and Theobalt, C. Neural sparse voxel fields. *Advances in Neural Information Processing Systems*, 33:15651–15663, 2020.
- Long, X., Liu, L., Theobalt, C., and Wang, W. Occlusion-aware depth estimation with adaptive normal constraints. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IX 16*, pp. 640–657. Springer, 2020.

-
- Long, X., Liu, L., Li, W., Theobalt, C., and Wang, W. Multi-view depth estimation using epipolar spatio-temporal networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8258–8267, 2021.
- Long, X., Lin, C., Wang, P., Komura, T., and Wang, W. Sparseneus: Fast generalizable neural surface reconstruction from sparse views. In *European Conference on Computer Vision*, pp. 210–227. Springer, 2022.
- Lu, T., Yu, M., Xu, L., Xiangli, Y., Wang, L., Lin, D., and Dai, B. Scaffold-gs: Structured 3d gaussians for view-adaptive rendering. *arXiv preprint arXiv:2312.00109*, 2023.
- Ma, Z., Teed, Z., and Deng, J. Multiview stereo with cascaded epipolar raft. In *European Conference on Computer Vision*, pp. 734–750. Springer, 2022.
- Mildenhall, B., Srinivasan, P., Tancik, M., Barron, J., Ramamoorthi, R., and Ng, R. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, 2020.
- Morgenstern, W., Barthel, F., Hilsmann, A., and Eisert, P. Compact 3d scene representation via self-organizing gaussian grids. *arXiv preprint arXiv:2312.13299*, 2023.
- Müller, T., Evans, A., Schied, C., and Keller, A. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4): 1–15, 2022.
- Navaneet, K., Meibodi, K. P., Koohpayegani, S. A., and Pirsavash, H. Compact3d: Compressing gaussian splat radiance field models with vector quantization. *arXiv preprint arXiv:2311.18159*, 2023.
- Poole, B., Jain, A., Barron, J. T., and Mildenhall, B. Dreamfusion: Text-to-3d using 2d diffusion. In *The Eleventh International Conference on Learning Representations*, 2022.
- Schönberger, J. L., Zheng, E., Frahm, J.-M., and Pollefeys, M. 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, pp. 501–518. Springer, 2016.
- Sun, P., Kretzschmar, H., Dotiwalla, X., Chouard, A., Patnaik, V., Tsui, P., Guo, J., Zhou, Y., Chai, Y., Caine, B., et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2446–2454, 2020.
- Tang, J., Ren, J., Zhou, H., Liu, Z., and Zeng, G. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023.
- Vakalopoulou, M., Chassagnon, G., Bus, N., Marini, R., Zacharaki, E. I., Revel, M.-P., and Paragios, N. Atlasnet: Multi-atlas non-linear deep networks for medical image segmentation. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2018: 21st International Conference, Granada, Spain, September 16–20, 2018, Proceedings, Part IV* 11, pp. 658–666. Springer, 2018.
- Wang, P., Liu, Y., Chen, Z., Liu, L., Liu, Z., Komura, T., Theobalt, C., and Wang, W. F2-nerf: Fast neural radiance field training with free camera trajectories. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4150–4159, 2023.
- Xu, Q. and Tao, W. Multi-scale geometric consistency guided multi-view stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5483–5492, 2019.
- Xu, Q., Xu, Z., Philip, J., Bi, S., Shu, Z., Sunkavalli, K., and Neumann, U. Point-nerf: Point-based neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5438–5448, 2022.
- Yan, Z., Low, W. F., Chen, Y., and Lee, G. H. Multi-scale 3d gaussian splatting for anti-aliased rendering. *arXiv preprint arXiv:2311.17089*, 2023.
- Yang, Z., Chen, Y., Wang, J., Manivasagam, S., Ma, W.-C., Yang, A. J., and Urtasun, R. Unisim: A neural closed-loop sensor simulator. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1389–1399, 2023.
- Yao, Y., Luo, Z., Li, S., Fang, T., and Quan, L. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 767–783, 2018.
- Yoo, J.-C. and Han, T. H. Fast normalized cross-correlation. *Circuits, systems and signal processing*, 28: 819–843, 2009.
- Yu, Z., Peng, S., Niemeyer, M., Sattler, T., and Geiger, A. Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. *Advances in neural information processing systems*, 35:25018–25032, 2022.
- Yu, Z., Chen, A., Huang, B., Sattler, T., and Geiger, A. Mip-splatting: Alias-free 3d gaussian splatting. *arXiv preprint arXiv:2311.16493*, 2023.

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

Zwicker, M., Pfister, H., Van Baar, J., and Gross, M. Ewa splatting. *IEEE Transactions on Visualization and Computer Graphics*, 8(3):223–238, 2002.