# GaRField++: Reinforced Gaussian Radiance Fields for Large-Scale 3D Scene Reconstruction

Hanyue Zhang[1], Zhiliu Yang[1, 2, *], Xinhe Zuo[1], Yuxin Tong[1], Ying Long[1], and Chen Liu[3]
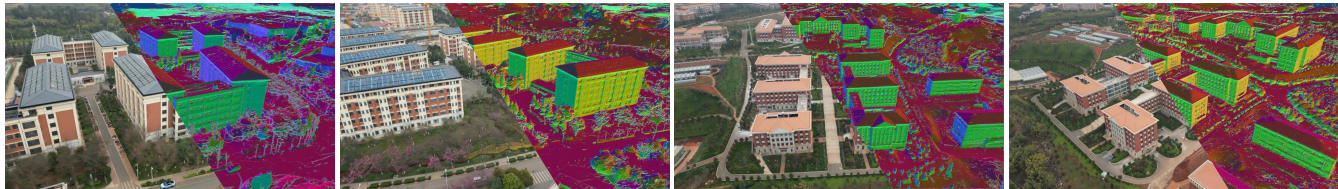
Fig. 1: Rendered RGB images and corresponding rendered depth normals from our GaRField++ framework on the self-collected data. Randomly rendered images from multiple views of the large-scale scenes are complete, smooth and detailed. This is achieved by constructing a divide-and-conquer Gaussian radiance field, which is reinforced by precisely modeling the color and opacity information and improving the training efficiency. The data is collected from the monocular camera of a DJI drone.

*Abstract*— This paper proposes a novel framework for large-scale scene reconstruction based on 3D Gaussian splatting (3DGS) and aims to address the scalability and accuracy challenges faced by existing methods. For tackling the scalability issue, we split the large scene into multiple cells, and the candidate point-cloud and camera views of each cell are correlated through a visibility-based camera selection and a progressive point-cloud extension. To reinforce the rendering quality, three highlighted improvements are made in comparison with vanilla 3DGS, which are a strategy of the ray-Gaussian intersection and the novel Gaussians density control for learning efficiency, an appearance decoupling module based on ConvKAN network to solve uneven lighting conditions in large-scale scenes, and a refined final loss with the color loss, the depth distortion loss, and the normal consistency loss. Finally, the seamless stitching procedure is executed to merge the individual Gaussian radiance field for novel view synthesis across different cells. Evaluation of Mill19, Urban3D, and MatrixCity datasets shows that our method consistently generates more high-fidelity rendering results than state-of-the-art methods of large-scale scene reconstruction. We further validate the generalizability of the proposed approach by rendering on self-collected video clips recorded by a commercial drone.

## I. INTRODUCTION

The recent advances in 3D reconstruction of large-scale urban scenes have reshaped modern society. It can serve as a visualization medium for AR/VR [1], aerial surveying [2], and city planning [3], [4], a high definition (HD) map for autonomous driving [5], [6], [7], [8], [9], [10], or a photorealistic simulator for unexpected cases in end-to-end autonomous driving and unmanned aerial vehicles (UAVs) [3], [11], [12], [13].

The task consists of high-fidelity reconstruction and real-time rendering for large areas that typically span more than

[1] School of Information Science and Engineering, Yunnan University, Kunming, Yunnan 650500, China.

[2] Yunnan Key Laboratory of Intelligent Systems and Computing, Yunnan University, Kunming, Yunnan 650500, China.

[3] Department of Electrical and Computer Engineering, Clarkson University, Potsdam, New York 13699, USA.

* Corresponding author, `zhiliu.yang@ynu.edu.cn`

1.5 $km^2$ [2]. In recent years, the field has been dominated by methods based on Neural Radiance Fields (NeRFs) [14]. Representative works include Block-NeRF [5], GgNeRF [15], Switch-NeRF [16] and Mega-NeRF [2]. However, these methods still lack the fidelity in preserving details. Recently, the 3D Gaussian Splatting (3DGS) technique [17] has gained significant attention for its outstanding performance in visual quality and rendering speed, achieving near-photorealistic rendering effects at 1080p resolution in real time. It has also been successfully applied to the reconstruction of dynamic scenes [18], [19], [20] and the generation of 3D content [21], [22]. However, the 3DGS still faces several scalability and accuracy challenges when dealing with large-scale environments.

Firstly, large-scale scenes typically encompass various objects, including the complex geometry structure such as grass, plants [23], and a large area of background such as the sky and water body [12]. Traditional 3DGS-based reconstruction methods do not adequately model normal depth and opacity information. Secondly, uneven lighting conditions in large-scale scenes may lead to significant appearance differences in captured images. When dealing with these variations, 3DGS tends to generate large-size 3D Gaussians with low opacity [12], which results in floating artifacts in novel views. Third, optimizing the entire large-scale scene requires multiple iterations, which become extremely time consuming and unstable without the proper regularization term and loss function design [24].

Recent efforts of large-scale scene reconstruction based on the 3DGS have mitigated some of the aforementioned shortcomings. Methods like visibility-based camera selection [12], appearance modeling [25], multimodal fusion [26], level of details [27] etc. are proposed correspondingly to improve the rendering quality. Although these methods produce reasonable results, they are still prone to some of the blurred area in the rendered views.

We propose a reinforced Gaussian radiance field for large-

scale 3D reconstruction, named GaRField++. We split the large-scale scene into multiple cells by following the Vast-Guaussian [12], then we implement visibility-based camera selection, relevant cameras from other cells and extended sets of the point cloud are enrolled for training to eliminate the floating artifacts. To enhance rendering fidelity, we leverage the ray-Gaussian-intersection volume rendering and improved density control strategies in the reconstruction of each cell. To mitigate uneven lighting conditions, we use a network architecture that integrates KAN [28] with convolutional neural networks (CNNs) to decouple appearance information. This color decoupling module is discarded after training to prevent impacting the rendering speed. In addition, a reinforced final loss is employed with color loss, depth distortion loss, and normal consistency loss.

In addition to testing on the challenging public dataset, we also utilize a DJI drone (Mini 3 Pro) to capture video clips from a large-scale scene to validate the effectiveness of our approach. Our contributions are as follows.

- GaRField++ is the first work to leverage the ray-Gaussian intersection volume rendering and the reinforced density control strategy for the large-scale 3D reconstruction, which consistently generates more high-fidelity rendering results than state-of-the-art methods.
- We leverage a color decoupling module based on KAN and CNN to address the appearance variations, enhancing the fidelity of the rendering results.
- We exploit the depth-normal consistency to construct the regulation term for large-scale area reconstruction, to increase continuity of 3DGS optimization.

## II. RELATED WORK

### A. Rendering with Radiance Fields

*1) Neural Radiance Fields:* Neural Radiance Fields (NeRF) [14] implicitly represents 3D scenes as a mapping from position and direction into radiance using a multi-layer perceptrons (MLPs), and achieves novel view synthesis through volumetric rendering techniques. Despite the significant progress made for 3D scene reconstruction and rendering by NeRF [14], they still face challenges in efficiency and memory usage when dealing with large-scale scenes. To improve rendering efficiency, researchers have proposed various strategies [29], [30], [31]. InstantNGP [29] firstly encodes the scene into a multi-resolution hashing table. Mip-NeRF [2] enhances NeRF's representation capacity for outdoor scenes by introducing the down sampling of conical frustums. Zip-NeRF [32] employs a hexagonal sampling strategy to address aliasing issues in the rendering.

*2) 3D Gaussian Splatting:* Rendering methods based on points utilize 3D Gaussian functions as geometric primitives, achieving the rapid rendering and a scene editing ability [17]. The 3D Gaussian Splatting (3DGS) further enhances rendering efficiency by employing optimized rasterization. Although 3DGS can produce high-fidelity 3D reconstruction results, methods such as Mip-splatting [33], LightGaussian [34], GSCore [35], Gaussianpro [36], Fregs [37], Eagles

[38], Compact3d [39] are proposed to improve the rendering process. Motivated by the method of EWA-Splatting [40], the Mip-Splatting [33] limits the frequency of the 3D representation and introduces a 2D Mip filter. Eagles [38], Compact3d [39], and others are committed to applying the VQ [41] trick to compress a large number of Gaussian primitives. Unlike FreGS [37], C3DGS [42], which optimizes on the software algorithms, GSCore [35] proposes a hardware acceleration unit to optimize the 3DGS pipeline in the rendering of the radiance field. GaussianPro [36] introduces an innovative paradigm for joint 2D-3D training to reduce the dependence on SfM initialization.

### B. Large-scale Scene Reconstruction

The neural rendering and the 3DGS-based rendering are naturally extended to the domain of large-scale scene reconstruction. Block-NeRF [5] divides large scenes into blocks and introduces appearance embeddings, learned pose refinement, and controllable exposure for the training of each individual block. Mega-NeRF [2] analyzes the data visibility of large-scale scenes, thereby proposes a sparse network structure where parameters are dedicated to different areas of the scene. Urban Radiance Fields [26] utilizes LiDAR and 2D optical flow data for large-scale scene reconstruction. Switch-NeRF [43] introduces a Mixture of Experts (MoE) system for end-to-end large-space modeling. A 3D point is assigned to an expert through a gating network, and the final rendering outcome is determined by the combined output of the expert and the gate value. VastGaussian [12] and CityGaussian [27] are representative works that take advantage of 3DGS for scalability and rendering fidelity of large-scale scene reconstruction. Additionally, DrivingGaussian [44] and StreetGaussians [25] aim at reconstructing large-scale dynamic scenes in autonomous driving using multi-modal data. StreetGaussians [25] uses Fourier transforms to effectively represent the temporal changes of spherical harmonics. DrivingGaussian [44] leverages the LiDAR priors and employs multi-frame multi-view data for hierarchical scene modeling. 3DGS-Calib [45] introduces LiDAR point clouds as reference points for Gaussian positions to construct a continuous scene representation.

While the aforementioned studies have effectively improved the rendering quality in the large-scale scene reconstruction compared to the methods proposed before inventing the NeRFs and the 3DGS, there is the space for improving the rendering precision of geometric structure and large homogeneous areas.

## III. METHODOLOGY

Our GaRField++ framework processes the input images through a structure-from-motions module, a scenes partitioning, a cells rendering, and a seamless stitching to construct a reinforced Gaussian radiance field, which gives its capability to synthesize photorealistic views. The overview of the entire framework is shown in Fig. 2.
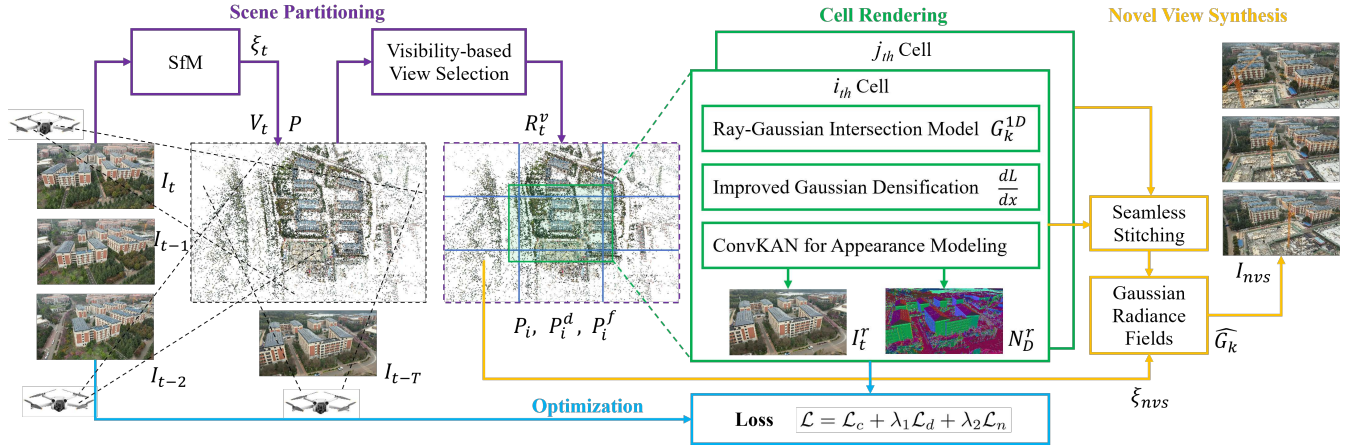
**Fig. 2: Overview of our GaRField++ framework. Scene Partitioning:** We implement a sparse reconstruction based on the Structure-from-Motion (SfM) method, generating a point cloud and estimating the initial camera pose for each image. Concurrently, we performed Manhattan alignment on the point cloud. Subsequently, we employ a coordinate-based regionalization and a visibility-based view selection strategy to split the point cloud. **Cell Rendering:** By leveraging the ray-Gaussian intersection model, enhanced Gaussian density control, and convolution KAN (Kernelized Attention Network)-based decoupled appearance modeling, we obtained the reconstruction results for each partition. **Optimization:** We employ a newly constructed loss function to optimize the training process. This loss function encompasses depth distortion loss, normal consistency loss, and color loss, thereby enhancing the accuracy and efficiency of large-scale reconstruction. **Novel View Synthesis:** we seamlessly stitched together the separate Gaussian fields from various cells to obtain a complete Gaussian field for the large-scale scene. This step enables the entire large-scale area model to support cross-border rendering, providing the possibility for the generation of novel view synthesis.

## A. Scenes Partitioning

We employ a divide-and-conquer strategy similar to [12] and [27], divide the large-scale scene into multiple cells, then render each cell independently.

*1) Sparse Reconstruction:* The input images of the scene are denoted as $\{I_t | t = 1, 2, ..., T\}$. Then the Structure-from-Motion (SfM) method, COLMAP [46], is adopted to generate a sparse point cloud $P$, and the initial camera pose $\xi_t$ is estimated for each image $I_t$. The camera views are defined as $V_t = \{I_t, \xi_t\}$. The Z axis of the point cloud $P$ is adjusted to be perpendicular to the ground plane by performing Manhattan world alignment [12].

*2) Visibility-based View Selection:* The best illumination condition and geometry visibility can be obtained by applying the coverage-wise point selection strategy, and details of the view selection are given below.

- **Coordinate-based Regionalization:** The large-scale scene is first divided into $N$ cells and we distribute parts of the point cloud to a specific cell. The point cloud within a cell is defined as $\{P_i | i = 1, 2, 3, ..., N\}$.
- **Point Clouds Extension:** Boundaries of the cell $i$ are expanded to enroll the common views between adjacent cells. The original bounded area of cell $i$ is $L_i^W \times L_i^H$, which now extends to $i$ is $(1+\beta)L_i^W \times (1+\beta)L_i^H$ by a certain percentage $\beta$. The set of point clouds $P_i$ is slightly dilated to $P_i^d$.
- **Cameras and Points Selection for Data Partitioning:** Given a cell $i$, the camera views from the adjacent cell $j$ is enrolled by checking the visibility criterion $R_t^v$, which is calculated by the following equation:

$$R_t^v = \{\frac{A_{proj}}{A_t^j} | A_t^j = W_t^j \times H_t^j\} \qquad (1)$$

Where $A_{proj}$ is the projected area of $i_{th}$ cell in image $I_t^j$ and $A_t$ is the area of pixels in image $I_t^j$ by multiplying the width of the image $W_t^j$ and height $H_t^j$. Cameras whose $R_t^v$ is larger than a predefined threshold [12] are selected to join the cell $i$. And more point cloud from the adjacent cell $j$ is selected in the current partition, only if those points can be observed from the newly added camera views $V_t$. The final point cloud inside a cell $i$ is further extended to $P_i^f$.

## B. Cells Rendering

The previous step produces the best point set, $P_i^f$, for modeling one of the partitions of large-scale areas, which represents a coarse description of the geometry distribution. Here, we further correlate these points with Gaussian primitives [17]. And our GaRField++ framework strengthens the radiance fields made up of Gaussian primitives with the following three reinforcements.

*1) Ray-Gaussian Intersection Model & Improved Gaussian Density Control:* The sparse point clouds of the scene is further depicted with a set of 3D Gaussian primitives $\{G_k | k = 1, \ldots, K\}$ correspondingly. The properties of each 3D Gaussian $G_k$ are parameterized by view-dependent color $\mathbf{c}_k \in \mathbb{R}^{3 \times 1}$, opacity $\alpha_k \in [0, 1]$, center $\mathbf{u}_k \in \mathbb{R}^{3 \times 1}$, scale $\mathbf{s} \in \mathbb{R}^{3 \times 1}$, and rotation $\mathbf{R} \in \mathbb{R}^{3 \times 3}$.

The Gaussian primitive $G_k$ of any point $\mathbf{x} \in \mathbb{R}^{3 \times 1}$ is depicted as:

$$G_k(\mathbf{x}) = \alpha_k e^{-\frac{1}{2}(\mathbf{x}-u_k)^T \Sigma_k^{-1}(\mathbf{x}-u_k)} \qquad (2)$$

Different from original 3DGS [17] method which projects Gaussian balls into 2D screen space and examine the Gaussian in 2D, ray-Gaussian intersection [23] is utilized here to convert 3D Gaussians at any point $\mathbf{x}$ into a 1D Gaussian $G_k^{1D}(\mathbf{x})$. For a given camera pose $\xi_t$ of the image $I_t$,

the contribution of Gaussian along its ray is defined as $\psi(G_k^{1D}, \xi_t)$. Then the color of a pixel $p_v$ in $I_t$ is rendered via alpha blending along the camera ray:

$$\mathbf{c}(p_v) = \sum_{k=1}^{K} \mathbf{c}_k \alpha_k \psi\left(\mathcal{G}_k^{1D}, \xi_t\right) \prod_{j=1}^{k-1} \left(1 - \alpha_j \psi\left(\mathcal{G}_j^{1D}, \xi_t\right)\right) \quad (3)$$

By utilizing the ray tracing volume rendering in Equation (3), the opacity along the ray is monotonically increasing until it reaches the maximal value.

Motivated by [23], an improved Gaussian densification strategy is used, in addition to the classical cloning or splitting, to handle areas that are overly blurred. To enlarge the gradients values, the magnitude of view position gradient is redesigned as:

$$\frac{dL}{d\mathbf{x}} = \sum_v \left\| \frac{dL}{dp_v} \frac{dp_v}{d\mathbf{x}} \right\| \quad (4)$$

where $\mathbf{x}$ is the center of Gaussian, $p_v$ is the pixels, and $\frac{dL}{d\mathbf{x}}$ is the position gradient of 3DGS [17]. Accumulating the norms $\| \cdot \|$ prevents the gradient signals from different pixels to negate each other. The densification strategy in our framework is executed at every certain iterations during the rendering.

*2) ConvKAN-based Decoupled Appearance Modeling:* To address the potential inconsistency between geometry and lighting in the rendering process, decoupled appearance modeling is required. The VastGaussian [12] utilizes a small CNN to predict the colors and illuminations of the images. Inspired by the Kernelized Attention Network (KAN) [28], [47], our decoupling network is designed by inserting KAN into CNNs. Replacing part of traditional convolution operations with KAN can improve rendering quality while keeping the model parameters almost unchanged.
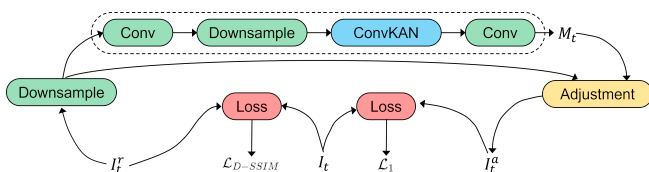


Fig. 3: Architecture of our ConvKAN-based decoupled appearance modeling.

As shown in Fig. 3, our decoupled appearance model consists of an initial convolutional layer that processes the initial input to extract preliminary features, a downsampling block, and a final convolutional layer where KAN replaces traditional convolution operations. The role of the downsampling block is to progressively downsample the feature maps, reducing the spatial resolution. The convKAN layer further processes the downsampled features and finally produces output through a sigmoid activation layer, with values ranging between 0 and 1. This color decoupling module is discarded after training, and thus it will not impact the rendering time.

*3) Optimization:* Rudimentary photo-metric loss is not reliable and effective for modeling large-scale reconstruction. Motivated by the regularization terms in 2DGS [24] and GOF [23], we optimize Gaussian model of $i_{th}$ cell with the following loss function:

$$\mathcal{L} = \mathcal{L}_c + \lambda_1 \mathcal{L}_d + \lambda_2 \mathcal{L}_n \quad (5)$$

$\mathcal{L}_d$ is the depth distortion loss proposed by 2DGS [24]. $\mathcal{L}_n$ is normal consistency loss, the normal $\mathbf{N}_D$ is estimated by the gradient of the depth map $D_t$. $\mathcal{L}_c$ is a RGB loss from 3DGS [17], which is defined as follow:

$$\mathcal{L}_c = \mathcal{L}_1\left(I_t^a, I_t\right) + \lambda_3 \mathcal{L}_{D-SSIM}\left(I_t^r, I_t\right) \quad (6)$$

As shown in Fig. 3, the $\mathcal{L}_{D-SSIM}$ metric predominantly penalizes deviations in structural integrity, and its application to the comparison between the rendered image $I_t^r$ and the original image $I_t$ ensures a high degree of appearance alignment between $I_t^a$ and $I_t$. Meanwhile, the task of recognizing appearance features is fulfilled by embeddings $L_t$ and our ConvKAN-based network. Furthermore, the loss function $\mathcal{L}_1$ is utilized to address the appearance discrepancies between the rendered image $I_t^a$ and the actual scene image $I_t$, accommodating ground truth images that may exhibit subtle variations in appearance compared to other images. After training, the rendered image $I_t^r$ is expected to maintain a consistent appearance with other images, enabling the Gaussian radiance field to learn the average appearance characteristics across all input views, as well as the precise geometric structure.

### C. Seamless Stitching & Novel view Synthesis

The Gaussian radiance fields within each cell is well-trained, and the Gaussian points outside the original boundary presented by $P_i$ (before the boundaries extension step) is cut out for seamless merging. Then we directly stitch different cells together, and the entire large-scale area model can now support cross-boarder rendering for novel view synthesis. Given a random camera pose $\xi_{nvs}$, the novel view $I_{nvs}$ can be rendered. To be noticed, $I_{nvs}$ are neither seen in the training datasets nor the testing dataset.

## IV. EXPERIMENTS

### A. Experimental Setup

*1) Dataset and Metrics:* The experiments are conducted across five large-scale scenarios: the Rubble and the Building from the Mill-19 dataset [2], the Residence from the Urban-Scene3D dataset [48], the small_city, which is a synthetic scene from the MatrixCity dataset [49], and Campus-YNU dataset collected by ourselves. The Campus-YNU dataset covers a region around $1\text{km} \times 1\text{km}$, which is captured simply using a DJI drone (Mini 3 Pro). SSIM, PSNR, and LPIPS [50] architecture serve as our evaluation metrics to quantitatively analyze the rendering results. All experiments were conducted on NVIDIA L40 GPUs with 48 GB memory for each card.
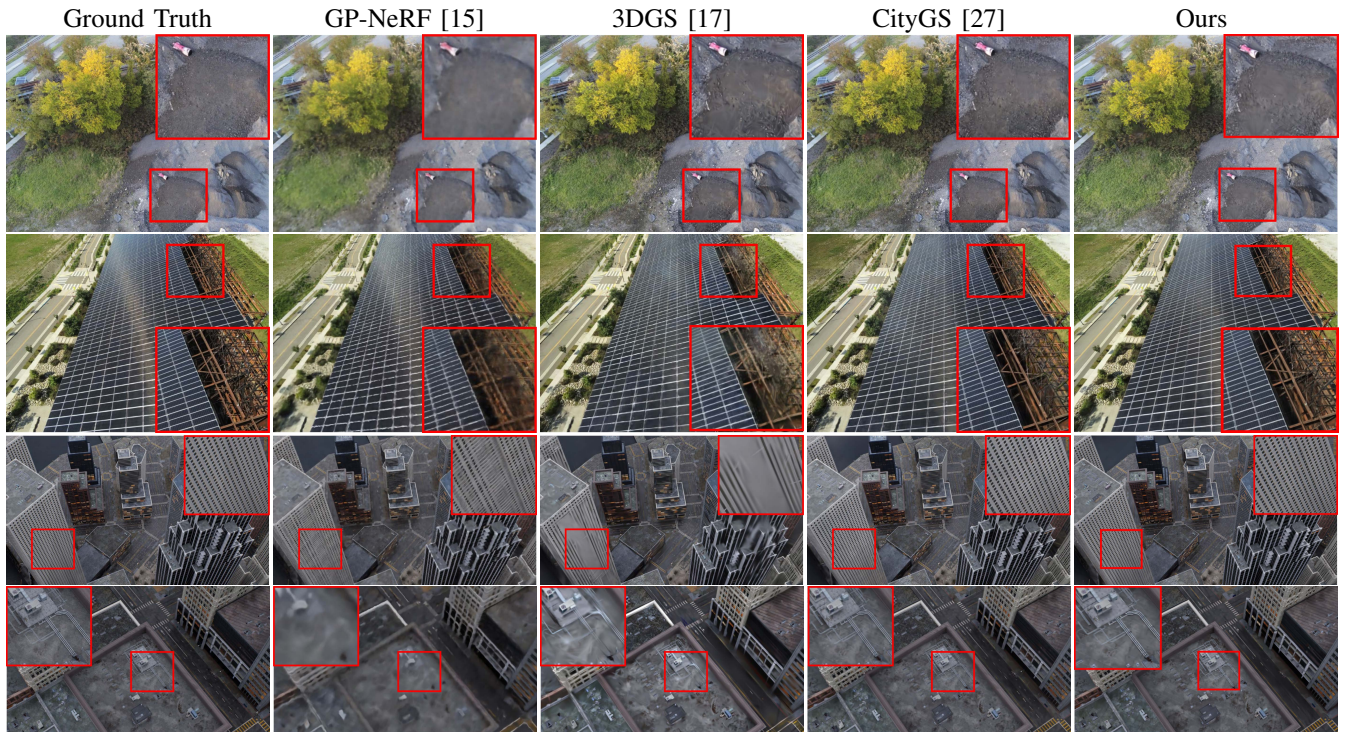
**Fig. 4: Qualitative Comparison with SOTA.** The first row represents the *Rubble* scenario, the second row manifests the *building* scenario, and the third and fourth rows showcase *small_city* scenes from the *MatrixCity* dataset. The experiment demonstrates superior capability of our GaRField++ framework in preserving color fidelity in rendered images, which is more closely resembling to the original images. Specifically, the region of interests are zoomed in with red box. **(Best viewed with zoom-in.)**

**TABLE I: Quantitative comparison on four challenging datasets with SOTA large-scale reconstruction methods.** Symbol '-' indicates that Mega-NeRF and Switch-NeRF are not evaluated on MatrixCity because of the difficulty in training them on different configuration. The red , orange and yellow colors respectively denote the best, the second best, and the third best results.

| Scenes | Building | | | Rubble | | | Residence | | | MatrixCity | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | SSIM↑ | PSNR↑ | LPIPS↓ | SSIM↑ | PSNR↑ | LPIPS↓ | SSIM↑ | PSNR↑ | LPIPS↓ | SSIM↑ | PSNR↑ | LPIPS↓ |
| Mega-NeRF [2] | 0.550 | 20.85 | 0.499 | 0.561 | 24.09 | 0.509 | 0.625 | 22.10 | 0.481 | - | - | - |
| Switch-NeRF[16] | 0.577 | 21.67 | 0.480 | 0.569 | 23.50 | 0.501 | 0.656 | 22.60 | 0.460 | - | - | - |
| GP-NeRF [15] | 0.570 | 21.10 | 0.489 | 0.565 | 24.32 | 0.489 | 0.659 | 22.29 | 0.450 | 0.610 | 23.60 | 0.392 |
| 3DGS [17] | 0.731 | 20.50 | 0.307 | 0.790 | 25.69 | 0.281 | 0.800 | 22.10 | 0.229 | 0.740 | 23.71 | 0.390 |
| CityGaussian [27] | 0.780 | 21.61 | 0.307 | 0.821 | 27.00 | 0.219 | 0.820 | 21.19 | 0.220 | 0.868 | 27.53 | 0.200 |
| Ours | **0.818** | **25.32** | **0.227** | **0.866** | **29.19** | **0.198** | **0.839** | **22.72** | **0.214** | **0.897** | **28.73** | **0.194** |

*2) Implementations and Baselines:* Our approach is compared to Mega-NeRF [2], Switch-NeRF [16], GP-NeRF [51], 3DGS [17], and CityGS [27]. First, we stop densification at 15k iterations as 3DGS [17]. Given that most of the data sets consist of a significantly larger number of input images than the data sets used in 3DGS [17], we adjusted the total number of training iterations to 60k. The default camera visibility is set to 0.25.

### B. Performance of Novel view Synthesis

*1) Comparison with SOTA:* As demonstrated in Table I, our method outperforms the state-of-the-art (SOTA) methods in terms of SSIM, PSNR, and LPIPS metrics for all four scenes (*Building*, *Rubble*, *Residence*, and *MatrixCity*. Here, *small_city* is selected from the *MatrixCity* dataset.). The qualitative results presented in Fig. 4 also validate the high fidelity of our rendering results. As shown in Fig. 4, our

rendering results achieve more realistic results, which is much closer to Ground Truth in the aspects of lighting and color. Specifically, in the building scenario, our renderings better preserve the detail of sunlight reflection on solar panels, which validates the effectiveness of our ray-Gaussian-intersection rendering, density control strategy, and the color decoupling module based on KAN and CNN.

*2) Experiments on Self-collected Data:* To validate the effectiveness and generalization capability of our framework in large-scale scenarios, we employ a DJI drone (Mini 3 Pro) to fly over a $1km \times 1km$ area, captured a dataset comprising 1,600 images at a resolution of $3768 \times 2118$ pixels. This scene was selected in our experiment for its intricate details, including solar panels, window arrays, and construction sites. Comparative experiments are conducted between 3DGS and our GaRField++. As shown in Table

II, the results of the experiment demonstrate an obvious improvement of our method over 3DGS considering the rendering quality. In Fig. 5, we can observe that our method makes solar panels and grasslands more precise compared to 3DGS.

**TABLE II: Experiments on Self-collected Campus-YNU Dataset.** We validate the effectiveness of our method compared to vanilla 3DGS on our self-collected data.

| Metrics | SSIM | PSNR | LPIPS |
|---------|------|------|-------|
| 3DGS [17] | 0.831 | 28.95 | 0.241 |
| Ours | **0.896** | **31.58** | **0.151** |



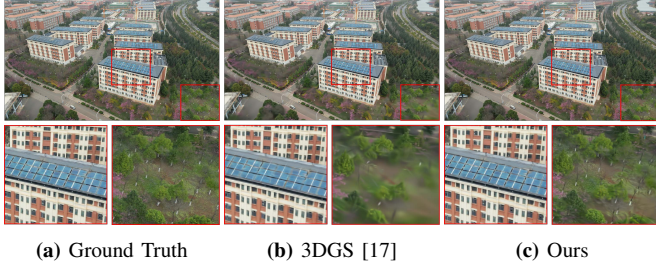**(a)** Ground Truth     **(b)** 3DGS [17]     **(c)** Ours

**Fig. 5: Comparison of our method with 3DGS on Self-collected data.** Fig. 5.a corresponds to the original image obtained from *Campus-YNU* scenario Fig. 5.b illustrates the image rendered using 3DGS, where the solar panels and the trees exhibit a degree of blurriness. Fig. 5.c demonstrates the image rendered with our proposed method, showing a decent enhancement in the clarity of the solar panels and the trees. **(Best viewed with zoom-in.)**

### C. Ablation Study

We conduct the ablation study on the Rubble scenario and our self-collected data to evaluate the different proposed techniques of our method. We randomly select 95% of the images from the Rubble dataset as the training set and 5% as testing set, and the ratio is kept same to our Campus-YNU dataset.

**TABLE III: Quantitative Results of the Ablation Study.** The red, orange and yellow colors respectively denote the best, the second best, and the third best results. Full stands for the configuration with CNN + KAN + ViS R2 + Full Loss.

| Scene | Rubble | | | Campus-YNU (Self-collect.) | | |
|-------|--------|--------|--------|--------|--------|--------|
| Metrics | SSIM | PNSR | LPIPS | SSIM | PNSR | LPIPS |
| Vis R0 | 0.853 | 29.18 | 0.592 | 0.889 | 30.56 | 0.161 |
| Vis R1 | 0.857 | **29.33** | 0.229 | 0.895 | 30.90 | 0.152 |
| CNN | 0.835 | 28.81 | 0.226 | 0.885 | 29.10 | 0.178 |
| $\mathcal{L}_c$ Only | 0.795 | 27.29 | 0.205 | 0.887 | 30.70 | 0.164 |
| Full | **0.854** | 29.09 | **0.201** | **0.896** | **31.58** | **0.151** |

*1) Camera Visibility Calculation:* As shown in the Table III, we conduct an investigation into camera visibility within a *Rubble* scenario. Specifically, we establish three levels of camera visibility, designated as Vis R0, Vis R1, and Vis R2, with settings of 0, 0.50, and 0.25, respectively. Throughout the experimental process, we use a full loss function and a color decoupling module integrated with CNN and KAN [28]. Subsequently, our approach employs the proposed camera visibility technique to test all these levels of visibility. As demonstrated in Fig. 6.d, Fig. 6.e, and Fig. 6.f, the camera visibility is found to be helpful in enhancing the quality of the rendering.

*2) Loss:* Our loss function, which is composed of depth distortion loss, normal consistency loss, and RGB loss derived from 3DGS [17], can better enhance the rendering quality of images compared to using only the RGB loss from 3DGS [17]. As illustrated in Fig. 6.h and Fig. 6.i, the rendered text is obviously clearer after using our loss function.

*3) Decoupled Color Model:* Our color decoupling module, which employs a network combining KAN [28] and CNN, has achieved superior results in reducing color variations in rendered images. As illustrated in Fig. 6.b and Fig. 6.c, compared to a color decoupling module composed solely of CNN, our approach can more effectively learn consistent geometric shapes and colors from training images with varying appearances. During the ablation experiments of the color decoupling module, we utilized a camera visibility of 0.25 and employed a full loss function.
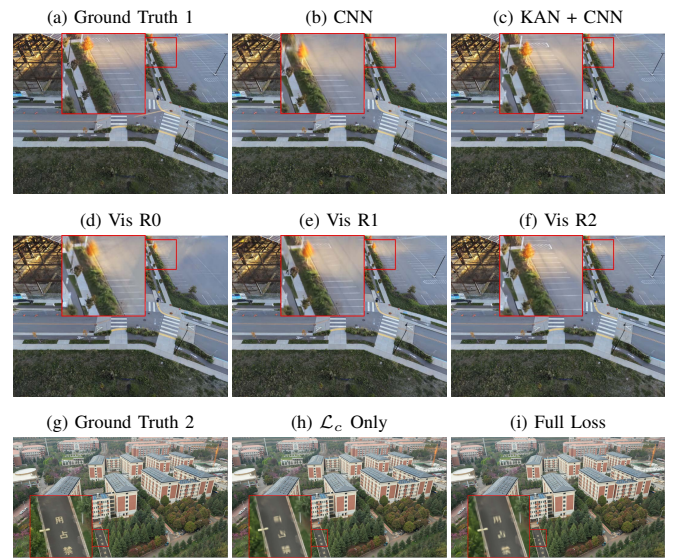


(a) Ground Truth 1    (b) CNN    (c) KAN + CNN

(d) Vis R0    (e) Vis R1    (f) Vis R2

(g) Ground Truth 2    (h) $\mathcal{L}_c$ Only    (i) Full Loss

**Fig. 6: Qualitative Results of the Ablation Study.** Ground Truth 1 represents the original image from the *Rubble* scenario. Ground Truth 2 corresponds to the original image obtained from the *Campus-YNU* scenario. **(Best viewed with zoom-in.)**

## V. CONCLUSION

In this work, we introduce GaRField++, a high-fidelity reconstruction and rendering method for large-scale scenes based on 3D Gaussian splatting. We employ a ray-Gaussian-intersection volume rendering and a density control strategy for large-scale reconstruction, a color decoupling module that combines KAN and CNN, a data partitioning method based on coordinates and camera visibility, and depth-normal consistency. We have achieved state-of-the-art rendering fidelity in mainstream benchmark tests and excellent rendering fidelity in our self-collected data set. However, we have not yet explored the optimal solutions for camera visibility and coordinate partitioning. In some scenarios, we still require hyper-parameter tuning to provide better rendering quality, and our model relies on the accuracy of the initial sparse point cloud. Additionally, our research may be applied to the 3D mesh extraction in the large-scale scenes. These works are left for our future endeavors.

## References

[1] L. Xu, V. Agrawal, W. Laney, T. Garcia, A. Bansal, C. Kim, S. Rota Bulò, L. Porzi, P. Kontschieder, A. Božič *et al.*, "Vr-nerf: High-fidelity virtualized walkable spaces," in *SIGGRAPH Asia 2023 Conference Papers*, 2023, pp. 1–12.

[2] H. Turki, D. Ramanan, and M. Satyanarayanan, "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*, 2022, pp. 12 922–12 931.

[3] H. Zhou, J. Shao, L. Xu, D. Bai, W. Qiu, B. Liu, Y. Wang, A. Geiger, and Y. Liao, "Hugs: Holistic urban 3d scene understanding via gaussian splatting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 21 336–21 345.

[4] Y. Bao, T. Ding, J. Huo, Y. Liu, Y. Li, W. Li, Y. Gao, and J. Luo, "3d gaussian splatting: Survey, technologies, challenges, and opportunities," *arXiv preprint arXiv:2407.17418*, 2024.

[5] M. Tancik, V. Casser, X. Yan, S. Pradhan, B. Mildenhall, P. P. Srinivasan, J. T. Barron, and H. Kretzschmar, "Block-nerf: Scalable large scene neural view synthesis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 8248–8258.

[6] C. Yan, D. Qu, D. Xu, B. Zhao, Z. Wang, D. Wang, and X. Li, "Gs-slam: Dense visual slam with 3d gaussian splatting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 19 595–19 604.

[7] H. Huang, L. Li, H. Cheng, and S.-K. Yeung, "Photo-slam: Real-time simultaneous localization and photorealistic mapping for monocular stereo and rgb-d cameras," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 21 584–21 593.

[8] N. Keetha, J. Karhade, K. M. Jatavallabhula, G. Yang, S. Scherer, D. Ramanan, and J. Luiten, "Splatam: Splat track & map 3d gaussians for dense rgb-d slam," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 21 357–21 366.

[9] V. Yugay, Y. Li, T. Gevers, and M. R. Oswald, "Gaussian-slam: Photo-realistic dense slam with gaussian splatting," *arXiv preprint arXiv:2312.10070*, 2023.

[10] M. Li, S. Liu, and H. Zhou, "Sgs-slam: Semantic gaussian splatting for neural dense slam," *arXiv preprint arXiv:2402.03246*, 2024.

[11] Z. Li, L. Li, and J. Zhu, "Read: Large-scale neural scene rendering for autonomous driving," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 2, 2023, pp. 1522–1529.

[12] J. Lin, Z. Li, X. Tang, J. Liu, S. Liu, J. Liu, Y. Lu, X. Wu, S. Xu, Y. Yan *et al.*, "Vastgaussian: Vast 3d gaussians for large scene reconstruction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 5166–5175.

[13] J. Gu, M. Jiang, H. Li, X. Lu, G. Zhu, S. A. A. Shah, L. Zhang, and M. Bennamoun, "Ue4-nerf: Neural radiance field for real-time rendering of large scene," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[14] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[15] L. Xu, Y. Xiangli, S. Peng, X. Pan, N. Zhao, C. Theobalt, B. Dai, and D. Lin, "Grid-guided neural radiance fields for large urban scenes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 8296–8306.

[16] M. Zhenxing and D. Xu, "Switch-nerf: Learning scene decomposition with mixture of experts for large-scale neural radiance fields," in *The Eleventh International Conference on Learning Representations*, 2022.

[17] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering." *ACM Trans. Graph.*, vol. 42, no. 4, pp. 139–1, 2023.

[18] J. Luiten, G. Kopanas, B. Leibe, and D. Ramanan, "Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis," *arXiv preprint arXiv:2308.09713*, 2023.

[19] Z. Yang, X. Gao, W. Zhou, S. Jiao, Y. Zhang, and X. Jin, "Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 20 331–20 341.

[20] G. Wu, T. Yi, J. Fang, L. Xie, X. Zhang, W. Wei, W. Liu, Q. Tian, and X. Wang, "4d gaussian splatting for real-time dynamic scene rendering," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 20 310–20 320.

[21] T. Yi, J. Fang, G. Wu, L. Xie, X. Zhang, W. Liu, Q. Tian, and X. Wang, "Gaussiandreamer: Fast generation from text to 3d gaussian splatting with point cloud priors," *arXiv preprint arXiv:2310.08529*, 2023.

[22] Z. Chen, F. Wang, Y. Wang, and H. Liu, "Text-to-3d using gaussian splatting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 21 401–21 412.

[23] Z. Yu, T. Sattler, and A. Geiger, "Gaussian opacity fields: Efficient and compact surface reconstruction in unbounded scenes," *arXiv preprint arXiv:2404.10772*, 2024.

[24] B. Huang, Z. Yu, A. Chen, A. Geiger, and S. Gao, "2d gaussian splatting for geometrically accurate radiance fields," in *ACM SIGGRAPH 2024 Conference Papers*, 2024, pp. 1–11.

[25] Y. Yan, H. Lin, C. Zhou, W. Wang, H. Sun, K. Zhan, X. Lang, X. Zhou, and S. Peng, "Street gaussians for modeling dynamic urban scenes," *arXiv preprint arXiv:2401.01339*, 2024.

[26] K. Rematas, A. Liu, P. P. Srinivasan, J. T. Barron, A. Tagliasacchi, T. Funkhouser, and V. Ferrari, "Urban radiance fields," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 932–12 942.

[27] Y. Liu, H. Guan, C. Luo, L. Fan, J. Peng, and Z. Zhang, "Citygaussian: Real-time high-quality large-scale scene rendering with gaussians," *arXiv preprint arXiv:2404.01133*, 2024.

[28] Z. Liu, Y. Wang, S. Vaidya, F. Ruehle, J. Halverson, M. Soljačić, T. Y. Hou, and M. Tegmark, "Kan: Kolmogorov-arnold networks," *arXiv preprint arXiv:2404.19756*, 2024.

[29] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *ACM transactions on graphics (TOG)*, vol. 41, no. 4, pp. 1–15, 2022.

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

[31] B. Deng, J. T. Barron, and P. P. Srinivasan, "JaxNeRF: an efficient JAX implementation of NeRF," 2020. [Online]. Available: https://github.com/google-research/google-research/tree/master/jaxnerf

[32] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, and P. Hedman, "Zip-nerf: Anti-aliased grid-based neural radiance fields," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 19 697–19 705.

[33] Z. Yu, A. Chen, B. Huang, T. Sattler, and A. Geiger, "Mip-splatting: Alias-free 3d gaussian splatting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 19 447–19 456.

[34] Z. Fan, K. Wang, K. Wen, Z. Zhu, D. Xu, and Z. Wang, "Lightgaussian: Unbounded 3d gaussian compression with 15x reduction and 200+ fps," 2023.

[35] J. Lee, S. Lee, J. Lee, J. Park, and J. Sim, "Gscore: Efficient radiance field rendering via architectural support for 3d gaussian splatting," in *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*, 2024, pp. 497–511.

[36] K. Cheng, X. Long, K. Yang, Y. Yao, W. Yin, Y. Ma, W. Wang, and X. Chen, "Gaussianpro: 3d gaussian splatting with progressive propagation," in *Forty-first International Conference on Machine Learning*, 2024.

[37] J. Zhang, F. Zhan, M. Xu, S. Lu, and E. Xing, "Fregs: 3d gaussian splatting with progressive frequency regularization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 21 424–21 433.

[38] S. Girish, K. Gupta, and A. Shrivastava, "Eagles: Efficient accelerated 3d gaussians with lightweight encodings," *arXiv preprint arXiv:2312.04564*, 2023.

[39] K. Navaneet, K. P. Meibodi, S. A. Koohpayegani, and H. Pirsiavash, "Compact3d: Compressing gaussian splat radiance field models with vector quantization," *arXiv preprint arXiv:2311.18159*, 2023.

[40] M. Zwicker, H. Pfister, J. Van Baar, and M. Gross, "Ewa volume splatting," in *Proceedings Visualization, 2001. VIS'01.* IEEE, 2001, pp. 29–538.

[41] W. H. Equitz, "A new vector quantization clustering algorithm," *IEEE transactions on acoustics, speech, and signal processing*, vol. 37, no. 10, pp. 1568–1575, 1989.

[42] S. Niedermayr, J. Stumpfegger, and R. Westermann, "Compressed 3d gaussian splatting for accelerated novel view synthesis," 2023.

[43] Z. Mi and D. Xu, "Switch-nerf: Learning scene decomposition with mixture of experts for large-scale neural radiance fields," in

*International Conference on Learning Representations (ICLR)*, 2023. [Online]. Available: https://openreview.net/forum?id=PQ2zoIZqvm

[44] X. Zhou, Z. Lin, X. Shan, Y. Wang, D. Sun, and M.-H. Yang, "Drivinggaussian: Composite gaussian splatting for surrounding dynamic autonomous driving scenes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 21 634–21 643.

[45] Q. Herau, M. Bennehar, A. Moreau, N. Piasco, L. Roldao, D. Tsishkou, C. Migniot, P. Vasseur, and C. Demonceaux, "3dgs-calib: 3d gaussian splatting for multimodal spatiotemporal calibration," *arXiv preprint arXiv:2403.11577*, 2024.

[46] J. L. Schonberger and J.-M. Frahm, "Structure-from-motion revisited," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4104–4113.

[47] A. D. Bodner, A. S. Tepsich, J. N. Spolski, and S. Pourteau, "Convolutional kolmogorov-arnold networks," *arXiv preprint arXiv:2406.13155*, 2024.

[48] Y. Liu, F. Xue, and H. Huang, "Urbanscene3d: A large scale urban scene dataset and simulator," 2021.

[49] Y. Li, L. Jiang, L. Xu, Y. Xiangli, Z. Wang, D. Lin, and B. Dai, "Matrixcity: A large-scale city dataset for city-scale neural rendering and beyond," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 3205–3215.

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

[51] Y. Zhang, G. Chen, and S. Cui, "Efficient large-scale scene representation with a hybrid of high-resolution grid and plane features," *arXiv preprint arXiv:2303.03003*, 2023.