

# Robust and Efficient 3D Gaussian Splatting for Urban Scene Reconstruction

Zhensheng Yuan<sup>1,2</sup> Haozhi Huang<sup>1</sup> Zhen Xiong<sup>1</sup> Di Wang<sup>1\*</sup> Guanghua Yang<sup>1\*</sup>

<sup>1</sup>Jinan University <sup>2</sup>University of Macau

zhensheng@stu2022.jnu.edu.cn, hzhuang@jnu.edu.cn, acxz2000@stu2022.jnu.edu.cn,  
diwang@jnu.edu.cn, ghyang@jnu.edu.cn

## Abstract

*We present a framework that enables fast reconstruction and real-time rendering of urban-scale scenes while maintaining robustness against appearance variations across multi-view captures. Our approach begins with scene partitioning for parallel training, employing a visibility-based image selection strategy to optimize training efficiency. A controllable level-of-detail (LOD) strategy explicitly regulates Gaussian density under a user-defined budget, enabling efficient training and rendering while maintaining high visual fidelity. The appearance transformation module mitigates the negative effects of appearance inconsistencies across images while enabling flexible adjustments. Additionally, we utilize enhancement modules, such as depth regularization, scale regularization, and anti-aliasing, to improve reconstruction fidelity. Experimental results demonstrate that our method effectively reconstructs urban-scale scenes and outperforms previous approaches in both efficiency and quality. The source code is available at: <https://yzslab.github.io/REUrbanGS>.*

## 1. Introduction

Urban-scale scene reconstruction plays a crucial role in various fields, such as autonomous driving [13, 27], urban planning [25], and digital twins [6], making it a highly active research area. The recently emerged 3D Gaussian Splatting (3DGS) [8] represents scenes with explicit 3D elliptical Gaussians, enabling high-fidelity, photorealistic reconstruction from image collections while supporting real-time rendering.

Compared to traditional MVS methods [37] relying on depth estimation and mesh reconstruction, 3DGS optimizes explicit 3D Gaussians for higher rendering fidelity. Unlike implicit neural representations [22, 27] that require costly volume rendering, 3DGS leverages efficient rasterization for fast reconstruction and real-time rendering. However, its

explicit representation introduces scalability challenges, as spatial complexity increases with scene size. For instance, reconstructing a 360° unbounded outdoor scene, like *Bicycle* from MipNeRF360 [1] requires over 6 million Gaussians, with out-of-memory (OOM) errors occurring beyond 11 million Gaussians on a 24GB GPU. While larger memory can temporarily mitigate OOM issues, excessive Gaussians burden rasterizers and demand longer training, making real-time rendering impractical.

Urban-scale datasets often exhibit greater complexity with limited controlled acquisition conditions [3, 21, 41]. Variations in temporal factors (e.g. time of day, seasonal changes), weather patterns, and sensor configurations induce significant appearance discrepancies for identical objects. In addition, transient elements such as pedestrians and vehicles are frequently captured. 3DGS [8] tends to overfit these localized variations by introducing superfluous Gaussian primitives that align solely with a very limited number training viewpoints. However, these redundant Gaussians manifest as floating artifacts and structural inconsistencies when rendered from novel perspectives, severely compromising reconstruction fidelity.

To address these challenges, we propose a novel, efficient, and robust 3DGS method specifically designed for urban scene reconstruction. Building upon the original 3DGS framework, our approach incorporates multiple enhancements at different stages to achieve efficient, high-quality reconstruction and real-time rendering.

**Promoting reconstruction efficiency:** We extend Block-NeRF [41] with a novel visibility-based image selection mechanism to minimize redundancy in data preprocessing (Section 3.2). Additionally, the densification strategy is refined to concentrate computational resources within individual partitions, thereby avoiding redundant computation in irrelevant regions and significantly improving efficiency (Section 3.3). To achieve real-time rendering under limited resources, we introduce a novel bottom-up LOD generation strategy. Each level is carefully controlled within a predefined budget and builds upon the previous one, effectively managing resource consumption during training. This strat-

\*Corresponding author

egy significantly improves reconstruction efficiency while maintaining superior quality compared to previous pruning or compression-based methods. The generated levels are dynamically selected during rendering using an LOD switching strategy, enabling real-time performance with limited resources (Section 3.4).

**Improving reconstruction quality:** An appearance transform module is designed to learn and mitigate appearance variations of objects across different images, effectively neutralizing their negative impact on reconstruction. Once the reconstruction is complete, this module allows for flexible transformations of the scene’s appearance without negatively impacting rendering speed. Furthermore, we also propose scale and depth regularization to mitigate the generation of floaters and artifacts to further improving the reconstruction quality (Section 3.5).

Experimental results demonstrate that our method outperform existing methods in terms of reconstruction quality, resource efficiency, and rendering speed, enabling the reconstruction of arbitrarily large urban-scale scenes. The main contributions are summarized as follows:

- We propose a novel visibility-based data division strategy and in-partition prioritized densification method, to achieve efficient urban-scale scene reconstruction.
- A controllable LOD generation strategy is designed to real-time dynamic LOD selection and real-time rendering under limited resources.
- A fine-grained appearance transform module is developed to allow flexibly adjust the appearance without affecting the real-time performance, significantly enhancing the robustness to inter-image appearance variations.

## 2. Related Works

### 2.1. Novel View Synthesis

**Neural Radiance Field (NeRF).** NeRF [22] uses neural networks as an implicit representation of scenes to enable photorealistic novel view synthesis, has achieved remarkable success and garnered significant attention from researchers. Subsequent works have sought to enhance its quality [1, 2, 21, 39], extent its application to dynamic scene [28–30, 38], surface reconstruction [14, 43, 49], and so on.

However, the rendering speed is a major obstacle to the broader adoption of NeRF. Rendering a single 1080P image takes several minutes, which falls far short of real-time requirements. Subsequent methods [7, 33, 40, 52] have improved rendering efficiency. Notably, Instant-NGP [24] introduced multi-resolution hash encoding, achieving orders-of-magnitude improvements in efficiency. Real-world scenes typically feature richer content and intricate details, but these NeRF-based methods, constrained by model capacity, often learn only coarse, low-frequency information.

**3D Gaussian Splatting.** 3DGS [8] has recently emerged as a groundbreaking method. It explicitly represents scenes using 3D elliptical Gaussians and efficiently rasterizes them into images, enabling photorealistic reconstruction and real-time rendering. It significantly surpasses NeRF-based approaches in rendering efficiency and visual quality, attracting considerable research interest. Numerous derivative works emerged that enhanced its reconstruction fidelity [10, 12, 31, 35, 53], extended it to dynamic scene [15, 19, 44, 48], and explored its robustness under diverse conditions [5, 35, 54]. Nonetheless, explicit representations inherently suffer from high computational resource demands for large-scale scenes. To address this, various methods [4, 11, 23, 26] have been proposed to compress the models. Recently, Taming3DGS [20] introduced a steerable densification strategy, enabling the control of 3DGS memory usage during the training phase while minimizing the negative impact in fidelity. Despite these advancements, the current focus of these methods remains on small-scale scenes, scaling 3DGS to truly urban-scale environments remains challenging due to rapidly increasing computational and memory demands.

### 2.2. Large Scale Scene Reconstruction

Several NeRF-based methods have proposed approaches for reconstructing large-scale scenes. Block-NeRF [41] and Mega-NeRF [42] apply a manual divide-and-conquer strategy, segmenting scenes into multiple regions represented by separate MLPs. Switch-NeRF [57] replaces manual partitioning with a learnable gating network. BungeeNeRF [45] progressively reconstructs urban scenes at multiple detail levels by dynamically adding modules. Grid-NeRF [46] employs a dual-branch structure with coarse-to-fine refinement based on multi-resolution hash encoding. GP-NeRF [55] combines 3D hash-grids and 2D plane features for higher accuracy and efficiency. However, these methods continue to face significant challenges in rendering speed and reconstruction fidelity, limiting their practical applicability in large-scale urban environments.

Several recent works have extended 3DGS to large-scale scenes. Grendel-GS [56] addresses resource constraints by leveraging multiple GPUs, yet scaling hardware with scene complexity remains impractical. VastGaussian [16] adopts a divide-and-conquer approach combined with CNN-based color transformations to handle appearance variations; however, it fails to address the real-time rendering challenge, and image-space transformations often cause unstable optimization and unnatural color artifacts. Hierarchical-3DGS [9] and CityGaussian [18] take a step further by introducing LOD strategies for real-time rendering of large-scale scenes. However, they lack mechanisms to control resources during training, relying instead on post-training compression and extensive fine-tuning, which inevitably

compromises quality in large-scale scenes due to higher compression requirements.

Different from previous methods, our approach ensures strict resource control during training, eliminating the need for post-training compression. Furthermore, our novel appearance transform module enables fine-grained adjustments at the Gaussian level, enhancing robustness and flexibility while maintaining real-time performance.

### 3. Method

#### 3.1. Preliminary

3DGS [8] represents a scene using a set of explicit 3D Gaussians. Each 3D Gaussian is defined by:

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

where  $\mu$  is the center and  $\Sigma$  is the covariance matrix that can be decomposed into a rotation  $\mathbf{R} \in SO(3)$  and a scale  $\mathbf{S} \in \mathbb{R}^3$ . Each Gaussian is assigned a color  $c$  and an opacity  $o$  to represent scene appearance. When rendering an image, the 3D Gaussians are first projected onto the image plane, forming 2D Gaussians  $G'(x)$ . These are then processed in depth order, to compute each pixel color using  $\alpha$ -blending:

$$C(x_p) = \sum_{i \in N} c_i o_i G'_i(x_p) \prod_{j=1}^{i-1} (1 - o_j G'_j(x_p)) \quad (2)$$

where  $x_p$  represents a pixel,  $o_i$  and  $c_i$  denote the opacity and color of the  $i$ -th Gaussian respectively, and  $N$  represents the set of Gaussians covering the pixel  $x_p$ .

During the training process, 3DGS periodically performs a densification operation on those Gaussians that exhibit relatively high mean gradients in the image space, thereby improving their ability to capture the underlying geometry. As described in [8], the loss function of 3DGS includes the L1 and D-SSIM metrics, computed between the rendered image  $\hat{I}$  and its corresponding ground-truth image  $I$ :

$$\mathcal{L} = (1 - \lambda)\mathcal{L}_1(I, \hat{I}) + \lambda\mathcal{L}_{D-SSIM}(I, \hat{I}) \quad (3)$$

where  $\lambda$  is a hyperparameter, which is set to 0.2 in [8].

By optimizing the attributes of the Gaussians and carrying out densification to minimize this loss, 3DGS ultimately fulfills its goal of reconstructing the target scene.

#### 3.2. Scene and Data Division

We partition the scene horizontally and then assign training images to them. Initially, images are assigned based on their locations: an image  $I_i$  is assigned to the  $j$ -th partition if its location falls within the partition's bounding box. For images located outside the bounding box, we calculate a point-based visibility with respect to the partition to determine whether they should be assigned to it. This strategy is illustrated in Figure 1.

##### 3.2.1. Point-based Visibility

Visibility is calculated using the 3D point cloud and its association with 2D feature points, both generated by Structure from Motion (SfM). For an unselected image  $I_i$ , the 3D point cloud of the scene is projected onto its image plane, and compute its convex hull area  $V_i$ . Subsequently, we leverage the relationship between the 2D feature points of the  $I_i$  and 3D points to extract those within the  $j$ -th partition. These extracted feature points are then used to calculate the convex hull area  $V_{ij}$ . Therefore, visibility is calculated by  $V_{ij}/V_i$ .

The feature points are inherently occlusion-aware, ensuring that only visible regions are considered. This effectively prevents redundant image selection, achieving higher quality with the same number of training iterations.

##### 3.2.2. Partition Rebalancing

In practice, the central partition typically contains more high-visibility images than the edge partitions, causing an unbalanced workload. To address this, we adjust the partition sizes after visibility-based data division: partitions with too few images are merged with the smallest neighboring partition, while those with too many images are subdivided. This process repeats iteratively until the image distribution is reasonably uniform.

#### 3.3. In-Partition Prioritized Densification

Excessive resource allocation to areas outside the partition is unnecessary during training. However, simply increasing the gradient threshold in these areas may lead to Gaussians within the partition shifting outward to compensate for under-reconstruction, thereby degrading the quality of partition boundaries. To solve this problem, as shown in Figure 2 we propose a distance-related threshold for each Gaussian:

$$\tau_i = \hat{\tau}_{\min} \left( \frac{\min(d_i, \hat{d}_{\max})}{\hat{d}_{\max}} \cdot (\eta - 1) + 1 \right) \quad (4)$$

Where the  $d_i$  is the distance between the  $i$ -th Gaussian and the partition,  $\hat{d}_{\max}$  is the distance at which the maximum threshold  $\hat{\tau}_{\max} = \hat{\tau}_{\min} \cdot \eta$  ( $\eta \geq 1$ ) is applied. The  $i$ -th Gaussian will only be considered for densification if and only if its mean gradient satisfies  $\bar{\Delta}_{G_i} > \tau_i$ .

This strategy effectively reduces resource consumption in out-of-partition regions, accelerating training while preserving final reconstruction quality.

#### 3.4. Controllable Level-of-detail

The original 3DGS densification strategy [8] lacks resource constraints, making it impractical for urban-scale scenes. To address this, we extend the steerable densification strategy of [20], generating multiple levels of detail in a bottom-up manner while strictly enforcing predefined resource lim-

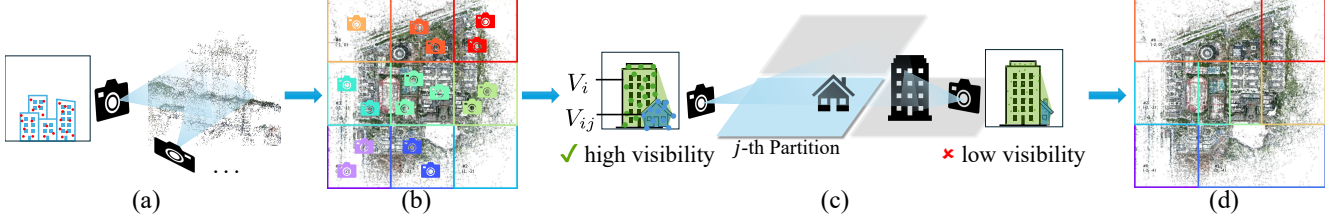


Figure 1. **The process of scene and data division.** (a) Obtain the 3D point cloud and its corresponding 2D feature points through estimating camera poses by SfM. (b) Determine training cameras based on their spatial locations after partitioning the scene into smaller regions. (c) Cameras outside the partitions are assigned based on visibility. The  $V_i$  and  $V_{ij}$  correspond to the green and blue regions in the image, respectively. Only cameras with high visibility are utilized for training. (d) Adjust the partition sizes to achieve a more balanced workload.

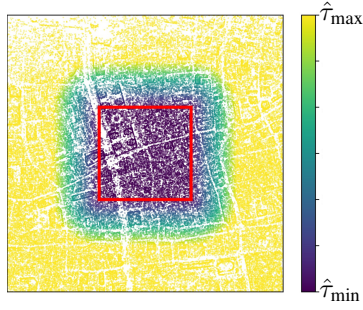


Figure 2. **In-partition prioritized densification.** The red rectangle is the partition bounding box, and points represent Gaussians. Point colors indicate gradient thresholds.

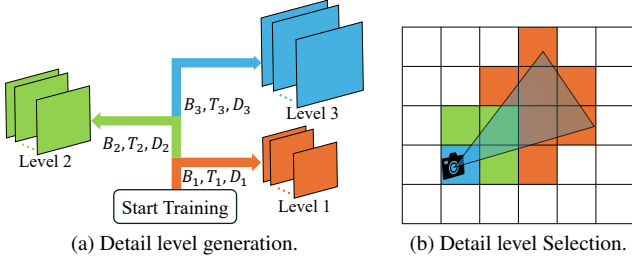


Figure 3. **Controllable LOD generation and detail level selection.** (a) During training, detail levels are progressively generated in a bottom-up manner, guided by resource budgets  $B$ , densification intervals  $T$ , and downsampling factors  $D$ . (b) During rendering, detail levels are dynamically selected based on the partition-camera distance, assigning higher levels to closer partitions and lower levels to distant ones. Invisible partitions are culled.

its. During rendering, appropriate levels are dynamically selected, ensuring efficient real-time rendering.

### 3.4.1. Controllable Detail Level Generation

The number of LOD levels  $l \in \mathbb{Z}^+$ , the budget  $B$ , densify interval  $T$  and image downsample factor  $D$  for each level are defined as:

$$\{B_i \in \mathbb{Z}^+ \mid B_1 < B_2 < \dots < B_l\} \quad (5)$$

$$\{T_i \in \mathbb{Z}^+ \mid T_1 > T_2 > \dots > T_l\} \quad (6)$$

$$\{D_i \in (0, 1] \mid D_1 < D_2 < \dots < D_l = 1\} \quad (7)$$

where the budget  $B$  is the parameter of the densification strategy of [20]. These parameters implies that when training at lower detail levels, a lower budget, longer densification intervals, and lower-resolution images will be utilized. This avoids unnecessary focus on high-frequency details when reconstructing lower levels.

As shown in Figure 3a, for each partition, the reconstruction begins with the 1st level. For the  $i$ -th level, upon completion of training, a checkpoint is created, and the budget, interval and downsample factor are changed to  $B_{i+1}$ ,  $T_{i+1}$  and  $D_{i+1}$  to facilitate the generation of the next level, until all levels are progressively generated.

This process is entirely end-to-end, eliminating the need for extensive post-processing steps common in compression strategy. Experiments show that our method achieves higher quality than compression-based method while enabling faster completion by utilizing low-resolution images and a smaller budget for lower levels.

### 3.4.2. Detail Level Selection

During rendering, we adopt the strategy proposed by City-Gaussian [18], performing detail level selection at the partition level, as illustrated in the Figure 3b. However, LOD selection incurs additional computational overhead. To further improve rendering efficiency, we also adopt the tile-based culling purposed by StopThePop [31] to disregard Gaussians with a low contribution in tiles.

## 3.5. Quality Enhancements

To further improve the reconstruction quality, a series of methods are introduced to overcome the key limitations of 3DGS in this subsection. Section 3.5.1 proposes the appearance transform module to ensure robust adaptation to appearance variations in images. To mitigate local overfitting and eliminate artifacts, Section 3.5.2 and Section 3.5.3 propose two regularization techniques. Section 3.5.4 adopts an anti-aliasing technique and enhances the fidelity of fine



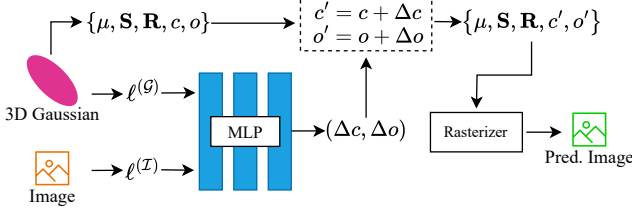


Figure 4. **Illustration of the appearance transform.** For each image and 3D Gaussian,  $\ell^{(G)}$  represents the Gaussian embedding and  $\ell^{(I)}$  represents the image embedding, respectively. By the lightweight MLP,  $\Delta c$  and  $\Delta o$  can be predicted to adjust the color  $c$  and opacity  $o$  of the 3D Gaussian.

details. Section 3.5.5 presents a method to remove transient objects.

### 3.5.1. Appearance Transform Module

3DGS [8] tends to generate floaters to overfit appearance variations across training images, degrading reconstruction quality. Existing methods like NeRF-W [21] and SWAG [5] mitigate appearance inconsistencies by assigning embeddings only to images, limiting flexibility and leading to suboptimal results. We propose a fine-grained appearance transform module that assigns embeddings to both individual images and each 3D Gaussian independently. As shown in Figure 4, these embeddings are processed by a lightweight MLP to predict per-Gaussian color and opacity offsets, enabling precise and adaptable appearance adjustments. This approach enhances robustness against inconsistencies, eliminates floating artifacts, and supports post-reconstruction appearance editing. Once computed, the offsets are reusable for subsequent frames as long as the image embedding remains unchanged, ensuring no additional rendering cost.

**Similarity regularization:** Adjacent Gaussians typically show similar variations. To leverage this property and avoid overfitting, we introduce a similarity regularization to encourage more similar embeddings among neighboring Gaussians. Since normalized embeddings are used, cosine similarity is adopted to compute the loss:

$$\mathcal{L}_{i,j}^{\text{sim}} = w_{i,j} \left( 1 - \frac{\ell_i^{(G)} \cdot \ell_j^{(G)}}{\|\ell_i^{(G)}\| \|\ell_j^{(G)}\|} \right) \quad (8)$$

$$\mathcal{L}_{\text{sim}} = \frac{1}{|M| \binom{k}{2}} \sum_{i \in M} \sum_{\substack{j, l \in \text{knn}_{i,k} \\ j < l}} \mathcal{L}_{j,l}^{\text{sim}} \quad (9)$$

where  $\ell_i^{(G)}$  and  $\ell_j^{(G)}$  are the embeddings of the  $i$ -th and  $j$ -th Gaussian, respectively.  $M$  is the set of Gaussians sampled randomly for regularization,  $\text{knn}_{i,k}$  denotes  $K$ -nearest neighbors of  $i$ ,  $w_{i,j}$  is the decay factor [19] and determined

by the distance between the two Gaussians:

$$w_{i,j} = \exp(-\lambda_w \|\mu_i - \mu_j\|) \quad (10)$$

where  $\lambda_w$  controls the decay rate, chosen based on the scale of the scene.

**Opacity offset regularization:** Based on real-world experience, most appearance transformations do not involve changes in transparency. To prevent the model from unnecessarily overusing transparency to fit color variations, we introduce an additional regularization term for the opacity offset, restricting transparency changes to only a few Gaussian components:

$$\mathcal{L}_{\Delta o} = \frac{1}{N} \sum_{i=1}^N \Delta o_i. \quad (11)$$

### 3.5.2. Scale Regularization

During optimization, we frequently observed scale-anomalous Gaussians, such as those exceeding the scene size or forming highly anisotropic shapes. These anomalies cause severe artifacts during camera rotation due to ordering inversion [31]. To address this, we introduce a scale regularization with two components: a maximum constraint to prevent excessive growth and a ratio constraint to maintain reasonable proportions. This regularization effectively suppresses artifacts caused by scale anomalies and improves the stability and consistency of the rendering results.

The maximum constraint limits the upper bound of Gaussian scales to prevent them from growing to unreasonable values:

$$\mathcal{L}_{\text{ms}} = \frac{\sum_i \mathbb{1}\{\mathbf{S}_i > s_{\text{max}}\} \cdot \mathbf{S}_i}{\sum_i \mathbb{1}\{\mathbf{S}_i > s_{\text{max}}\} + \delta} \quad (12)$$

where the  $\mathbb{1}$  is an indicator function that takes the value 1 when the condition is true and 0 otherwise,  $\mathbf{S}_i$  is the scales of the  $i$ -th Gaussian,  $s_{\text{max}}$  is the maximum acceptable scale,  $\delta$  prevents division by zero.

The ratio constraint enforces a limit on the ratio between the first and second largest scales of a Gaussian, avoiding highly anisotropic shapes:

$$r_i = \frac{\max(\mathbf{S}_i)}{\text{median}(\mathbf{S}_i)} \quad (13)$$

$$\mathcal{L}_r = \frac{\sum_i \mathbb{1}\{r_i > r_{\text{max}}\} \cdot r_i}{\sum_i \mathbb{1}\{r_i > r_{\text{max}}\} + \delta} \quad (14)$$

where  $r_{\text{max}}$  is the maximum acceptable ratio.

### 3.5.3. Depth Regularization

Inspired by the DNGaussian [12], we utilize Depth Anything V2 [47] to predict fine-grained depth maps from RGB images and align them to actual depths using the SfM point

cloud. During training, we alternate between [12]’s hard depth and soft depth regularization. The depth loss  $\mathcal{L}_d$  is the L1 error between the rendered and estimated depths. This effectively mitigates the floating artifacts, significantly enhancing visual realism.

### 3.5.4. Anti-aliasing and Detail Enhancement

To improve rendering quality, we integrate anti-aliasing from Mip-Splatting [53] and adopt AbsGS [51] to enhance fine details. These adaptations effectively reduce artifacts and boost visual fidelity.

### 3.5.5. Transient Objects Removal

In practical data collection, the presence of transient objects (e.g., pedestrians, vehicles) is unavoidable. These transient objects introduce visual artifacts into the reconstructed scene, adversely affecting the final reconstruction quality. We utilize an open-world object detection model [34] to identify the 2D bounding boxes of these transient objects and employ them as prompts for a semantic segmentation model [32] to generate fine-grained masks.

## 3.6. Loss of Individual Partition Training

The loss for partition training consists of five components:

$$\mathcal{L}' = \mathcal{L} + \lambda_{\text{sim}}\mathcal{L}_{\text{sim}} + \lambda_{\Delta o}\mathcal{L}_{\Delta o} + \lambda_d\mathcal{L}_d + \lambda_s(\mathcal{L}_{\text{ms}} + \mathcal{L}_r) \quad (15)$$

Where  $\mathcal{L}$  is Equation (3). The  $\lambda_{\text{sim}} = 0.2$ ,  $\lambda_{\Delta o} = 0.05$  and  $\lambda_s = 0.05$ . The initial value of  $\lambda_d$  is 0.5, with an exponential decay scheduler reduces it to a final value of 0.01.

## 4. Experiments

### 4.1. Experimental Setup

**Datasets.** Our method is evaluated on three aerial photography scenes, including the *Rubble* scene from Mega-NeRF [42] and two self-collected scenes, *JNU-ZH* and *BigCity*. Notably, we also conducted validation using *Building* scene from Mega-NeRF [42] as well as *Residences*, *Sci-Art* and *Campus* scenes from UrbanScene3D [17], with results provided in the supplementary materials E.2. The *JNU-ZH* was collected over three months using two different cameras, consisting of approximately 5,000 images and covering an area of about 300,000 m<sup>2</sup>. This scene exhibits diverse appearance variations due to seasonal changes, weather conditions, and lighting differences. The *BigCity* was captured using five cameras, comprising around 10,000 images and spanning approximately 5 km<sup>2</sup>. It features a complex road network and a dense distribution of diverse buildings. We use three detail levels for these scenes.

**Metrics.** We conduct quantitative comparisons of the rendered images using three metrics: PSNR, SSIM, and VGGNet-based LPIPS. To evaluate rendering performance, we measured the FPS and the average number of Gaussians

(denoted as #G, in 10<sup>6</sup>) required for rendering the test set at the same resolution, where #G directly reflects VRAM consumption. All methods were evaluated using a single NVIDIA A100-80G GPU.

## 4.2. Results

Our proposed method is compared against four existing methods: Switch-NeRF [57], CityGaussian [18], Hierarchical-3DGS [9], and 3DGS [8]. Quantitative results are summarized in Table 1. The first section of the table compares our method, with the LOD mode disabled, against other methods without LOD mode or with it disabled. For quality-related metrics (SSIM, PSNR, and LPIPS), the results indicate that our method outperforms others. The only exception is the *Rubble* scene, where the LPIPS score matches that of CityGaussian. The second section of the table evaluates our method with LOD mode enabled. Compared to other LOD-enabled methods, our method consistently outperforms previous approaches across all three quality-related metrics. Moreover, in most cases, the results with LOD mode enabled surpass the non-LOD results of other methods. This underscores our method’s ability to achieve high-fidelity reconstructions of urban-scale scenes.

Regarding efficiency-related metrics-#G and FPS, while the #G in our method is not the smallest, it remains within a reasonable range and supports real-time rendering within 24 GB of memory. Notably, #G can be further reduced by lowering the budget  $B$  of the LOD generation. Meanwhile, the FPS does not experience a significant decline and consistently ranks as either the best or second-best, making real-time rendering entirely feasible. By comparing the results of our method with and without LOD mode, it becomes evident that the number of Gaussians is significantly reduced, leading to a substantial increase in FPS. Meanwhile, the quality experiences only minimal degradation. This can be attributed to our bottom-up detail level generation strategy, allows better preservation of fine-grained geometric and texture details, which is especially beneficial for challenging urban-scale scenes. Notably, other methods exhibit a lower #G in the *BigCity* scene, mainly due to our adjustment of their hyperparameters to ensure execution within 80 GB memory. Further details are provided in Supplementary Material C. The visualization results are shown in Figure 5, demonstrating that our method achieves superior detail recovery and exhibits a greater ability to eliminate artifacts.

### 4.3. LOD Generation

We conducted experiments on *Rubble* to analyze the impact of different values of the budget  $B$  for detail level generation. Table 2 presents the results. As observed, reducing  $B$  decreases the number of Gaussians, thereby lowering resource consumption, but at the cost of reconstruction qual-

Scene	Rubble					JNU-ZH					BigCity				
Metrics	SSIM	PSNR	LPIPS	#G	FPS	SSIM	PSNR	LPIPS	#G	FPS	SSIM	PSNR	LPIPS	#G	FPS
Switch-NeRF	0.544	23.05	0.508	—	<0.1	0.574	21.96	0.587	—	<0.1	0.469	20.39	0.645	—	<0.1
CityGaussian (no LOD)	<u>0.813</u>	<u>25.77</u>	<b>0.228</b>	<u>9.60</u>	57.0	<u>0.776</u>	<u>22.57</u>	0.282	<u>11.76</u>	34.6	0.825	<u>24.57</u>	<u>0.240</u>	<u>58.47</u>	<u>23.1</u>
Hierarchical-3DGS (no LOD)	0.768	23.76	0.275	11.13	49.6	0.772	21.23	<u>0.260</u>	17.42	29.8	—	—	—	—	—
3DGS	0.796	25.72	0.304	<b>7.10</b>	<b>108.9</b>	0.763	22.02	0.350	<b>3.66</b>	<b>152.7</b>	<u>0.830</u>	24.52	0.279	<b>33.21</b>	22.7
Ours (no LOD)	<b>0.826</b>	<b>27.29</b>	<b>0.228</b>	13.52	<u>78.2</u>	<b>0.822</b>	<b>25.85</b>	<b>0.232</b>	25.58	<u>41.2</u>	<b>0.847</b>	<b>26.62</b>	<b>0.219</b>	75.15	<b>23.7</b>
CityGaussian	<u>0.785</u>	<u>24.90</u>	<u>0.256</u>	<b>2.95</b>	<b>105.2</b>	<u>0.770</u>	<u>22.33</u>	0.286	<b>3.27</b>	<b>69.2</b>	0.712	22.24	0.344	<b>3.04</b>	<b>122.5</b>
Hierarchical-3DGS	0.741	23.38	0.300	7.23	57.4	0.760	21.12	<u>0.274</u>	11.21	34.6	<u>0.775</u>	<u>23.17</u>	<u>0.289</u>	17.09	3.2
Ours	<b>0.814</b>	<b>27.03</b>	<b>0.245</b>	<u>3.60</u>	<u>99.7</u>	<b>0.816</b>	<b>25.71</b>	<b>0.240</b>	<u>6.65</u>	<u>63.9</u>	<b>0.838</b>	<b>26.41</b>	<b>0.231</b>	<u>6.84</u>	<u>73.0</u>

Table 1. **Quantitative results on three large scene datasets.** We report SSIM $\uparrow$ , PSNR $\uparrow$ , LPIPS $\downarrow$ , the number of Gaussians (#G, in  $10^6$ ) $\downarrow$  and FPS $\uparrow$  on test views. The **best** and second best results are highlighted. All missing results are denoted by a “—”.

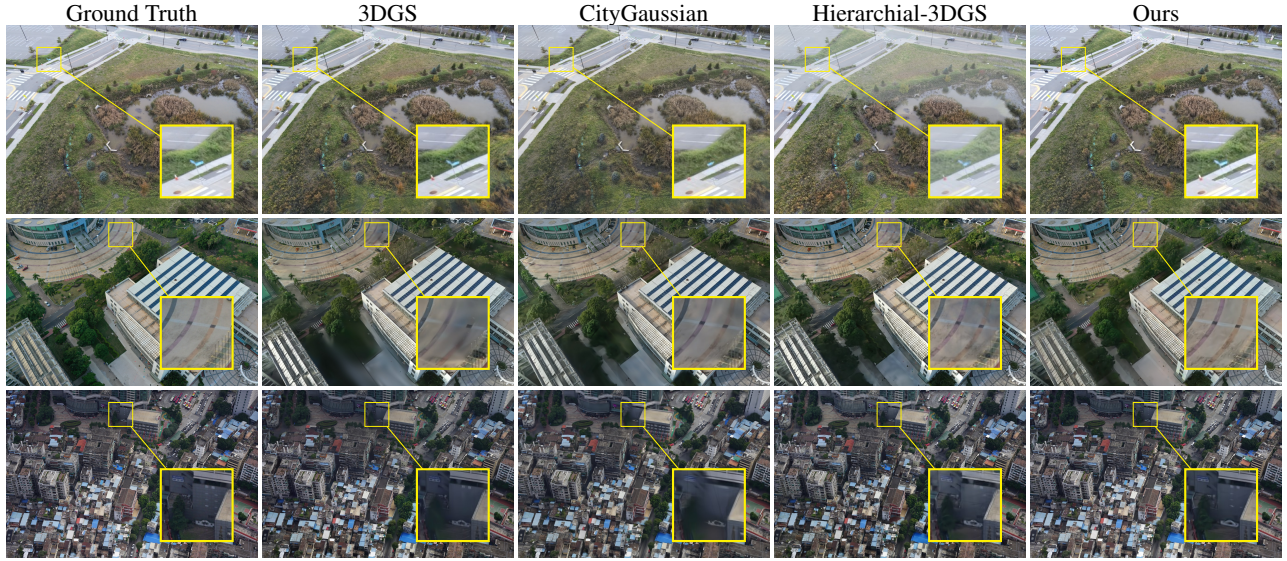


Figure 5. **Visualization results.** All methods (excluding 3DGS) render in LOD mode. Ours demonstrates better detail preservation and fewer artifacts.

ity. However, increasing  $B$  to beyond a certain threshold does not necessarily improve quality, because  $B$  only imposes an upper limit, and the scene may not require as many Gaussians as the upper bound allows.

Budget ( $\times 100$ )	SSIM	PSNR	LPIPS	#G	FPS
(1024, 2048, 4096)	0.771	26.13	0.297	<b>1.61</b>	<b>126.9</b>
(2048, 4096, 8192)	0.797	26.65	0.265	<u>2.69</u>	<u>108.8</u>
(4096, 8192, 16384)	<u>0.814</u>	<u>27.03</u>	<u>0.245</u>	3.60	99.7
(8192, 16384, 32768)	<b>0.816</b>	<b>27.11</b>	<b>0.242</b>	3.80	96.4

Table 2. **Quantitative evaluation of budget  $B$  for detail level generation.** Adjusting the budget effectively controls resource consumption, but also impacts the quality.

#### 4.4. Ablation Study

We conduct ablation experiments to evaluate the impact of different components of our proposed method.

**Point-based Visibility.** The 1st row of Table 3 presents the results without point-based visibility, meaning that cam-

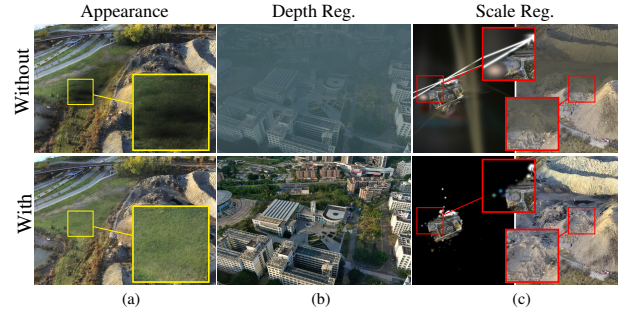


Figure 6. **Visualization results of ablation study.** Our proposed components effectively suppress the artifacts.

eras are assigned to partitions solely based on spatial locations. We follow the setup of Block-NeRF [41] and expand the bounding box of each partition by 50% to define the range for selecting training cameras based on their locations. However, it indicates that a purely position-based selection leads to a suboptimal result. As shown in Ta-



Scene	<i>Rubble</i>					<i>JNU-ZH</i>					<i>BigCity</i>				
Metrics	SSIM	PSNR	LPIPS	#G	FPS	SSIM	PSNR	LPIPS	#G	FPS	SSIM	PSNR	LPIPS	#G	FPS
w/o vis.	0.803	26.95	0.256	4.01	94.3	0.809	25.14	0.242	8.53	58.6	0.826	25.74	0.238	6.92	<u>73.7</u>
w/o appearance	0.771	25.17	0.284	4.65	82.6	0.780	22.57	0.246	8.30	54.8	0.831	25.08	0.233	6.89	<b>79.5</b>
w/o depth reg.	<b>0.817</b>	26.99	<b>0.240</b>	<b>3.46</b>	<u>98.2</u>	0.815	25.54	0.241	<b>6.56</b>	<u>63.7</u>	<b>0.845</b>	<u>26.40</u>	<b>0.221</b>	6.83	70.0
w/o scale reg.	0.814	26.93	0.241	3.65	89.4	<b>0.824</b>	<b>25.81</b>	<b>0.230</b>	7.07	56.5	0.829	25.42	0.235	6.85	68.4
w/o in-prior.	<b>0.817</b>	<b>27.22</b>	<b>0.240</b>	3.90	97.5	0.815	<u>25.72</u>	0.241	7.90	61.2	<u>0.840</u>	26.33	<u>0.228</u>	<b>6.80</b>	73.0
full	0.814	<u>27.03</u>	0.245	<u>3.60</u>	<b>99.7</b>	<u>0.816</u>	25.71	<u>0.240</u>	<u>6.65</u>	<b>63.9</b>	0.838	<b>26.41</b>	0.231	6.84	73.0

Table 3. **Qualitative ablations.** The results correspond to the removal of point-based visibility, appearance transform module, depth regularization, scale regularization, and in-partition prioritized densification, respectively.

Scene	<i>Rubble</i>	<i>JNU-ZH</i>	<i>BigCity</i>
w/o vis.	1.70×	1.46×	1.28×
full		1×	

Table 4. **Impact of the point-based visibility on the number of training cameras.** The component effectively reduces camera redundancy.

Scene	<i>Rubble</i>		<i>JNU-ZH</i>		<i>BigCity</i>	
Metrics	Time	#G	Time	#G	Time	#G
w/o in-prior.	1.49×	1.86×	1.64×	2.05×	1.30×	1.37×
full				1×		

Table 5. **Impact of in-partition prioritized densification on training time and the number of Gaussians.** It effectively reduces the number of Gaussians while accelerating training.

ble 4, our camera selection strategy can significantly reduce redundant cameras, thus achieving better results under the same number of training iterations.

**Appearance Transform Module.** The 2nd row of Table 3 demonstrates the effect of omitting the appearance transform model. This model significantly improves all three quality metrics across all scenes. Additionally, it sometimes noticeably reduces the number of Gaussians, as it prevents 3DGS from introducing extra Gaussians to overfit appearance variations. This, in turn, eliminates artifacts, as shown in Figure 6a.

**Depth Regularization.** The 3rd row of Table 3 presents the results without depth regularization. Although it does not have a significant impact on metrics, it remains a crucial component. Without it, the model tends to produce discrete floaters that overfit to certain training cameras. As shown in Figure 6b, these floaters often obstruct the scene when rendering from viewpoints that differ significantly from the training cameras, ultimately affecting the visual experience.

**Scale Regularization.** As shown in the 4th row of Table 3, scale regularization may lead to a slight decline in quality sometimes, because it prevents from generating scale-abnormal Gaussians that fit variations in appearance. But without it will lead to severe artifacts, as shown in Figure 6c.

**In-Partition Prioritized Densification.** The 5th row of Table 3 shows that the our densification strategy has only a minimal impact on overall quality, particularly with respect to the SSIM and LPIPS metrics, where enabling or disabling it yields only negligible differences. However, as illustrated in Table 5, it markedly accelerates the training process.

## 5. Conclusion

We propose a robust and efficient 3D Gaussian Splatting method tailored for urban-scale scene reconstruction. Our scene partitioning and visibility-based image selection enable scalable reconstruction within limited resources. The controllable LOD strategy provides precise resource regulation and real-time rendering. Additionally, our fine-grained appearance transform module and scale regularization significantly enhance visual consistency and flexibility. Extensive experiments demonstrate our method’s superior reconstruction quality, efficiency, and practical applicability to large-scale urban scenes.

Our method currently relies on accurate camera poses obtained from external SfM methods. Inaccurate or noisy poses can negatively impact reconstruction quality, especially in large-scale urban scenes. Enhancing robustness to pose inaccuracies is thus an important future direction. Additionally, our LOD switching mechanism currently is not incremental, increasing storage and computational overhead. Future work could explore incremental switching mechanisms for smoother transitions and improved resource efficiency.

**Acknowledgments.** This work was supported in part by the Science and Technology Development Fund, Macau SAR (Grants No. 0087/2022/AFJ and No. 001/2024/SKL), in part by the National Natural Science Foundation of China (Grant No. 62261160650), in part by the Research Committee of University of Macau (Grant No. MYRG-GRG2023-00116-FST-UMDF), and in part by the the Fundamental Research Funds for the Central Universities (Grant No. 21625360).



## References

- [1] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5470–5479, 2022. 1, 2
- [2] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Zip-nerf: Anti-aliased grid-based neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19697–19705, 2023. 2
- [3] Xingyu Chen, Qi Zhang, Xiaoyu Li, Yue Chen, Ying Feng, Xuan Wang, and Jue Wang. Hallucinated neural radiance fields in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12943–12952, 2022. 1
- [4] Yihang Chen, Qianyi Wu, Weiyao Lin, Mehrtash Harandi, and Jianfei Cai. Hac: Hash-grid assisted context for 3d gaussian splatting compression. In *European Conference on Computer Vision*, pages 422–438. Springer Nature Switzerland Cham, 2024. 2
- [5] Hiba Dahmani, Moussab Bennehar, Nathan Piasco, Luis Roldao, and Dzmitry Tsishkou. Swag: Splatting in the wild images with appearance-conditioned gaussians. In *European Conference on Computer Vision*, pages 325–340. Springer, 2024. 2, 5
- [6] Tianhu Deng, Keren Zhang, and Zuo-Jun Max Shen. A systematic review of a digital twin city: A new pattern of urban governance toward smart cities. *Journal of management science and engineering*, 6(2):125–134, 2021. 1
- [7] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5501–5510, 2022. 2
- [8] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 1, 2, 3, 5, 6
- [9] Bernhard Kerbl, Andreas Meuleman, Georgios Kopanas, Michael Wimmer, Alexandre Lanvin, and George Drettakis. A hierarchical 3d gaussian representation for real-time rendering of very large datasets. *ACM Transactions on Graphics (TOG)*, 43(4):1–15, 2024. 2, 6
- [10] Shakiba Kheradmand, Daniel Rebain, Gopal Sharma, Weiwei Sun, Yang-Che Tseng, Hossam Isack, Abhishek Kar, Andrea Tagliasacchi, and Kwang Moo Yi. 3d gaussian splatting as markov chain monte carlo. *Advances in Neural Information Processing Systems*, 37:80965–80986, 2025. 2
- [11] Marc Levoy and Pat Hanrahan. Light field rendering. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pages 441–452. 2023. 2
- [12] Jiahe Li, Jiawei Zhang, Xiao Bai, Jin Zheng, Xin Ning, Jun Zhou, and Lin Gu. Dngaussian: Optimizing sparse-view 3d gaussian radiance fields with global-local depth normalization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20775–20785, 2024. 2, 5, 6
- [13] Wei Li, CW Pan, Rong Zhang, JP Ren, YX Ma, Jin Fang, FL Yan, QC Geng, XY Huang, HJ Gong, et al. Aads: Augmented autonomous driving simulation using data-driven algorithms. *Science robotics*, 4(28):eaaw0863, 2019. 1
- [14] Zhaoshuo Li, Thomas Müller, Alex Evans, Russell H Taylor, Mathias Unberath, Ming-Yu Liu, and Chen-Hsuan Lin. Neuralangelo: High-fidelity neural surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8456–8465, 2023. 2
- [15] Zhan Li, Zhang Chen, Zhong Li, and Yi Xu. Spacetime gaussian feature splatting for real-time dynamic view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8508–8520, 2024. 2
- [16] Jiaqi Lin, Zhihao Li, Xiao Tang, Jianzhuang Liu, Shiyong Liu, Jiayue Liu, Yangdi Lu, Xiaofei Wu, Songcen Xu, Youliang Yan, et al. Vastgaussian: Vast 3d gaussians for large scene reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5166–5175, 2024. 2
- [17] Liqiang Lin, Yilin Liu, Yue Hu, Xingguang Yan, Ke Xie, and Hui Huang. Capturing, reconstructing, and simulating: the urbanscene3d dataset. In *European Conference on Computer Vision*, pages 93–109. Springer, 2022. 6, 2
- [18] Yang Liu, Chuanchen Luo, Lue Fan, Naiyan Wang, Junran Peng, and Zhaoxiang Zhang. Citygaussian: Real-time high-quality large-scale scene rendering with gaussians. In *European Conference on Computer Vision*, pages 265–282. Springer Nature Switzerland Cham, 2024. 2, 4, 6
- [19] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. In *2024 International Conference on 3D Vision (3DV)*, pages 800–809. IEEE, 2024. 2, 5
- [20] Saswat Subhajyoti Mallick, Rahul Goel, Bernhard Kerbl, Markus Steinberger, Francisco Vicente Carrasco, and Fernando De La Torre. Taming 3dgs: High-quality radiance fields with limited resources. In *SIGGRAPH Asia 2024 Conference Papers*, pages 1–11, 2024. 2, 3, 4
- [21] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7210–7219, 2021. 1, 2, 5
- [22] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 1, 2
- [23] Wieland Morgenstern, Florian Barthel, Anna Hilsmann, and Peter Eisert. Compact 3d scene representation via self-organizing gaussian grids. In *European Conference on Computer Vision*, pages 18–34. Springer Nature Switzerland Cham, 2024. 2
- [24] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a mul-

- tiresolution hash encoding. *ACM transactions on graphics (TOG)*, 41(4):1–15, 2022. [2](#)
- [25] Przemyslaw Musialski, Peter Wonka, Daniel G Aliaga, Michael Wimmer, Luc Van Gool, and Werner Purgathofer. A survey of urban reconstruction. In *Computer graphics forum*, pages 146–177. Wiley Online Library, 2013. [1](#)
- [26] Simon Niedermayr, Josef Stumpegger, and Rüdiger Westermann. Compressed 3d gaussian splatting for accelerated novel view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10349–10358, 2024. [2](#)
- [27] Julian Ost, Fahim Mannan, Nils Thuerey, Julian Knodt, and Felix Heide. Neural scene graphs for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2856–2865, 2021. [1](#)
- [28] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5865–5874, 2021. [2](#)
- [29] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *arXiv preprint arXiv:2106.13228*, 2021.
- [30] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10318–10327, 2021. [2](#)
- [31] Lukas Radl, Michael Steiner, Mathias Parger, Alexander Weinrauch, Bernhard Kerbl, and Markus Steinberger. Stopthepop: Sorted gaussian splatting for view-consistent real-time rendering. *ACM Transactions on Graphics (TOG)*, 43(4):1–17, 2024. [2](#), [4](#), [5](#)
- [32] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. [6](#)
- [33] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 14335–14345, 2021. [2](#)
- [34] Tianhe Ren, Qing Jiang, Shilong Liu, Zhaoyang Zeng, Wenlong Liu, Han Gao, Hongjie Huang, Zhengyu Ma, Xiaohe Jiang, Yihao Chen, et al. Grounding dino 1.5: Advance the “edge” of open-set object detection. *arXiv preprint arXiv:2405.10300*, 2024. [6](#)
- [35] Sara Sabour, Lily Goli, George Kopanas, Mark Matthews, Dmitry Lagun, Leonidas Guibas, Alec Jacobson, David J Fleet, and Andrea Tagliasacchi. Spotlessplats: Ignoring distractors in 3d gaussian splatting. *arXiv preprint arXiv:2406.20055*, 2024. [2](#)
- [36] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016. [1](#)
- [37] Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. 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*, pages 501–518. Springer, 2016. [1](#)
- [38] Liangchen Song, Anpei Chen, Zhong Li, Zhang Chen, Lele Chen, Junsong Yuan, Yi Xu, and Andreas Geiger. Nerf-player: A streamable dynamic scene representation with decomposed neural radiance fields. *IEEE Transactions on Visualization and Computer Graphics*, 29(5):2732–2742, 2023. [2](#)
- [39] Liang Song, Guangming Wang, Jiuming Liu, Zhenyang Fu, Yanzi Miao, et al. Sc-nerf: Self-correcting neural radiance field with sparse views. *arXiv preprint arXiv:2309.05028*, 2023. [2](#)
- [40] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5459–5469, 2022. [2](#)
- [41] Matthew Tancik, Vincent Casser, Xincheng Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretschmar. Block-nerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8248–8258, 2022. [1](#), [2](#), [7](#)
- [42] Haithem Turki, Deva Ramanan, and Mahadev 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*, pages 12922–12931, 2022. [2](#), [6](#)
- [43] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689*, 2021. [2](#)
- [44] Guanjuan Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang. 4d gaussian splatting for real-time dynamic scene rendering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20310–20320, 2024. [2](#)
- [45] Yuanbo Xiangli, Lining Xu, Xingang Pan, Nanxuan Zhao, Anyi Rao, Christian Theobalt, Bo Dai, and Dahua Lin. Bungeenerf: Progressive neural radiance field for extreme multi-scale scene rendering. In *European conference on computer vision*, pages 106–122. Springer, 2022. [2](#)
- [46] Lining Xu, Yuanbo Xiangli, Sida Peng, Xingang Pan, Nanxuan Zhao, Christian Theobalt, Bo Dai, and Dahua Lin. Grid-guided neural radiance fields for large urban scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8296–8306, 2023. [2](#)
- [47] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *Advances in Neural Information Processing Systems*, 37:21875–21911, 2025. [5](#)

- [48] Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20331–20341, 2024. [2](#)
- [49] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34:4805–4815, 2021. [2](#)
- [50] Vickie Ye, Ruilong Li, Justin Kerr, Matias Turkulainen, Brent Yi, Zhuoyang Pan, Otto Seiskari, Jianbo Ye, Jeffrey Hu, Matthew Tancik, et al. gsplat: An open-source library for gaussian splatting. *arXiv preprint arXiv:2409.06765*, 2024. [1](#)
- [51] Zongxin Ye, Wenyu Li, Sidun Liu, Peng Qiao, and Yong Dou. Absgs: Recovering fine details in 3d gaussian splatting. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 1053–1061, 2024. [6](#), [1](#)
- [52] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5752–5761, 2021. [2](#)
- [53] Zehao Yu, Anpei Chen, Binbin Huang, Torsten Sattler, and Andreas Geiger. Mip-splatting: Alias-free 3d gaussian splatting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19447–19456, 2024. [2](#), [6](#)
- [54] Dongbin Zhang, Chuming Wang, Weitao Wang, Peihao Li, Minghan Qin, and Haoqian Wang. Gaussian in the wild: 3d gaussian splatting for unconstrained image collections. In *European Conference on Computer Vision*, pages 341–359. Springer, 2024. [2](#)
- [55] Yuqi Zhang, Guanying Chen, and Shuguang Cui. Efficient large-scale scene representation with a hybrid of high-resolution grid and plane features. *Pattern Recognition*, 158: 111001, 2025. [2](#)
- [56] Hexu Zhao, Haoyang Weng, Daohan Lu, Ang Li, Jinyang Li, Aurojit Panda, and Saining Xie. On scaling up 3d gaussian splatting training. *arXiv preprint arXiv:2406.18533*, 2024. [2](#)
- [57] MI Zhenxing and Dan 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. [2](#), [6](#)

# Robust and Efficient 3D Gaussian Splatting for Urban Scene Reconstruction

## Supplementary Material

### A. Our Datasets

The *JNU-ZH* and *BigCity* scenes were collected by our team using drones, and their contents are shown in Figure 7. We employ COLMAP’s hierarchical SfM [36] to perform sparse reconstruction for both scenes. After finishing reconstruction, we use COLMAP’s geo-registration to align the reconstructed model with GPS coordinates. Subsequently, we compute the Euclidean distance between the estimated camera positions and their corresponding GPS coordinates. Outliers with excessively large distances are discarded, as they typically result from inaccurate pose estimations. This filtering process helps mitigate the negative impact of erroneous data. Finally, we downsample the images to a maximum edge length of 1600 pixels for experimentation. When partitioning, the sizes used for these two scenes are 180m and 400m, respectively.

### B. Implementation Details of Our Method

We implemented our method based on gsplat [50], which offers higher computational and memory efficiency compared to the [8].

The visibility threshold for dataset division is  $1/6$ . We use three detail levels for all scenes. The first and second levels each last a base of 15,000 iterations, with densification enabled. The third level runs for a base of 30,000 iterations, where densification is applied in the first half and the second half is solely dedicated to optimizing properties. Table 6 presents the hyperparameters for detail level generation in different scenes.

Scenes	$(B_1, B_2, B_3) \times 100$	$(T_1, T_2, T_3)$	$(D_1, D_2, D_3)$
<i>Rubble</i>	(4096, 8192, 16384)		(1/2, 1/3, 1)
<i>JNU-ZH</i>	(4096, 8192, 20480)	(300, 200, 100)	(1/4, 1/2, 1)
<i>BigCity</i>	(2048, 8192, 20480)		(1/4, 1/2, 1)

Table 6. The hyperparameters for the detail level generation.

In practice, these iteration counts and densification interval  $T$  are adjusted proportionally based on the number of images  $N$  in each partition, scaled by a factor of  $\max(N/600, 1)$ .

In the appearance transform model, the MLP consist of 1 hidden layers with 32 channels, followed by a ReLU activation. The output layer followed by a sigmoid activation. The Gaussian embeddings is 16-dimensional, while the image embedding is 32-dimensional. The initial learning rate for the MLP and embeddings is set to 0.01, and an exponential decay scheduler reduces it to a final value of 0.00025.

Every 50 iterations, we sample 20,480 Gaussians and select  $k = 16$  nearest neighbors to perform similarity regularization, minimizing the computational overhead.

In the scale regularization, the value of  $s_{\max}$  is set to a value corresponding to the typical size of most buildings in the scene, and  $r_{\max} = 10$ .

In the in-partition prioritized densification, The value of  $\hat{d}_{\max}$  is identical to the partition size. The maximum gradient threshold factor is  $\eta = 4$ . The minimum threshold is  $\tau_{\min} = 0.0002$  for the 1st and 2nd levels, and is 0.6 for the 3rd level with AbsGS [51] enabled.

### C. Hyperparameters of Other Methods

For the 3DGS, large-scale scenes generally require more iterations for sufficient optimization. Therefore, the training process was extended to 50 epochs, with densification enabled during the first 25 epochs. We also set the densification interval to  $1/6$  of an epoch, ensuring a consistent number of densifications across all scenes. When the number of input images is 600, these adjustments yield consistent hyperparameter with the original settings.

For Switch-NeRF, we utilized the official open-source implementation with its provided hyperparameters. When conducting experiments on our own scenes, we proportionally increased the number of training iterations based on the number of input images.

For the remaining methods, we utilized their official open-source implementations and use a similar number of partitions to reconstruct all the scenes. When evaluating the LOD mode of Hierarchical-3DGS, we used a granularity setting value of 6 pixels.

Due to the large scale and intricate details of the *BigCity* scene, none of the previous 3DGS-based methods can complete the experiment within an 80GB memory limit. Therefore, we made additional adjustments to the hyperparameters for these methods. For 3DGS, we double the densification gradient threshold. For CityGaussian, during the coarse training stage, we tripled the densification cycle and doubled the densification gradient threshold compared to the original settings. During the pruning stage, we increased the pruning ratios from the default 40%, 50%, and 60% to 70%, 80%, and 90%. For Hierarchical-3DGS, the excessive number of Gaussians made it impossible to evaluate the non-LOD mode. When evaluating its metrics under the LOD mode, we doubled the granularity settings from 6 pixels to 12 pixels. In contrast, our method can complete all steps, except for the non-LOD mode, with memory usage not exceeding 24GB.





Figure 7. **Our datasets: *JNU-ZH* and *BigCity*.**

## D. Appearance Transform Module

### D.1. Metric Calculation

Given that we have the appearance transform model, which optimizes only the embeddings of training set images, we followed a strategy similar to NeRF-W [21] to evaluate the test set images. Specifically, when computing the metrics for test images, we first optimize the image embedding  $\ell^{(\mathcal{I})}$  using the left half of the image and compute the metrics using the right half. Each partition is transformed using the embedding of the test image optimized within that partition, ensuring appearance consistency. Then, we optimize the embedding from scratch using the right half and computed the metrics with the left half. Finally, the average of the results from both rounds was taken as the final metric value for the entire image. This approach prevents information leakage and ensures fairness in the evaluation process. In practice, we further smooth transitions between partitions via weighted averaging.

### D.2. Appearance Transformation

After reconstruction, our method enables scene appearance transformation. Using the image embedding  $\ell^{(\mathcal{I})}$  of a training image, we can synthesize novel views that match its appearance. As shown in Figure 8, this enables transforming between different states of a building in the *JNU-ZH* scene.

## E. Additional Experiments

### E.1. Training Time Comparison

Table 7 presents the training times of all 3DGS-based methods. Except for 3DGS, all results were obtained under parallel training setups. The results show that our method is also competitive in terms of training efficiency, consistently

ranking among the best or second-best. We are not consistently the fastest due to the additional overhead introduced by the Appearance Transform Module, anti-aliasing, and various regularization mechanisms. Nonetheless, maintaining such competitive efficiency despite these added components demonstrates the effectiveness of our optimization strategies.

It is worth noting that the *Rubble* and *JNU-ZH* are relatively small, where parallel training provides limited benefits. In contrast, the *BigCity* is significantly larger, and the parallel setup leads to a substantial speedup.

Scenes	<i>Rubble</i>	<i>JNU-ZH</i>	<i>BigCity</i>
3DGS	1.45	3.21	67.39
CityGaussian	2.33	2.49	4.01
Hierarchical-3DGS	<b>1.00</b>	<b>2.18</b>	<b>1.71</b>
Ours	<u>1.30</u>	<b>2.14</b>	<u>2.01</u>

Table 7. **Comparison of training time (in hours).** Except for 3DGS, the results of all other methods were obtained under parallel training mode. VastGaussian is not included as it is not open-sourced.

### E.2. Additional Quantitative Comparison

Table 8 presents experimental results for the *Building* [42], *Residence*, *Sci-Art* and *Campus* [17] scenes. The camera poses are provided by Mega-NeRF. Overall, our method demonstrates a clear advantage in nearly all quality-related metrics. Although our method is not the most optimal in terms of resource consumption and rendering speed, it remains within a reasonable range and is close to the best-performing approach. It fully ensures real-time rendering. It is worth noting that our method can further reduce resource consumption by lowering the budget  $B$ . Figure 9

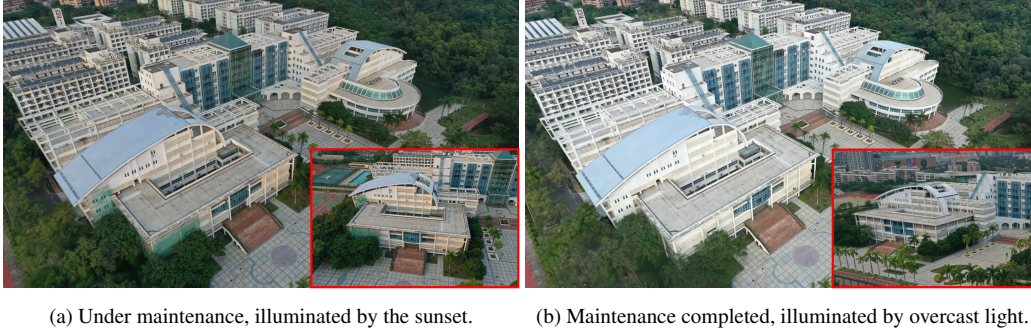


Figure 8. **Synthesis the two corresponding states from a new viewpoint based on the embedding vector provided by the reference image (bottom right).**

presents the visualization results for both scenes, demonstrating that our method achieves higher detail preservation and fewer artifacts.

### E.3. Quantitative Comparison of Detail Levels

Table 9 compares the performance of our three detail levels. The results show that all three levels achieve high reconstruction quality. The lower levels exhibit higher numerical values than the higher levels because they are trained and evaluated using downsampled images. Additionally, comparing the #G across levels further confirms that our LOD strategy effectively controls resource consumption.

### E.4. Additional ablations

Table 10 presents the results of ablation studies on the anti-aliasing, AbsGS, and tile-based culling components in our method.

**Anti-aliasing.** The 1th row of Table 10 reports the impact of anti-aliasing techniques. As shown in Figure 10a, it effectively prevents jagged edges from appearing in areas with low detail levels in images. However, this comes at the cost of requiring more Gaussians. In the *BigCity* scene, since the number of Gaussians has already reached the upper limit without anti-aliasing, enabling anti-aliasing does not allow for additional Gaussians, leading to a slight degradation in metrics. Nevertheless, this feature remains beneficial as it significantly enhances the visual experience.

**AbsGS.** The 2nd row of Table 10 provides the results obtained without AbsGS, highlighting its contribution to quality metrics and enhanced detail restoration, as shown in Figure 10b. While it does increase the number of Gaussians in certain scenarios, the increase remains within an acceptable range, ensuring that real-time rendering can still be achieved within the constraints of 24GB of VRAM.

**Tile-based culling.** The 3rd row of Table 10 presents the results without tile-based culling, illustrating its role in rendering efficiency. Tile-based culling noticeable improves rendering speed without negatively impacting ren-

dering quality. This is because it skips over redundant Gaussians with minimal contribution, ensuring efficiency while maintaining the desired quality.

### E.5. Comparison of Gaussian Embedding Lengths

We evaluate the impact of the length of Gaussian embedding on the *JNU-ZH* scene. As shown in Table 11, reducing the length leads to a slight drop in metrics, but it is acceptable if the goal is to reduce memory consumption.

### E.6. Comparison of LOD Selection Parameters

We evaluate the impact of our LOD parameters on the *JNU-ZH* scene. Table 12 presents the impact of different rendering-time partition sizes on metrics. It can be observed that while smaller partition sizes effectively reduce the number of Gaussians and lower resource consumption, they also lead to a certain degree of metric degradation. Additionally, a larger number of partitions incurs higher overhead due to the LOD selection. In contrast, larger partition sizes exhibit the opposite behavior. This suggests that an optimal partition size must strike a balance between efficiency and rendering quality. Table 13 illustrates the impact of different distance thresholds for detail levels. Increasing the distance thresholds generally improves rendering quality but also leads to higher resource consumption and reduced rendering speed. Therefore, selecting an appropriate distance threshold requires a trade-off between efficiency and quality.

### E.7. Evaluation of Similarity Regularization

Similarity regularization enhances the generalization ability of the appearance transformation module. When performing out-of-domain inference, such as predicting the unobserved regions of an image using its embedding, this regularization effectively mitigates abrupt color changes and suppresses artifacts, as illustrated in the first row of Figure 11.



Scene	Building					Residence					Sci-Art					Campus				
Metrics	SSIM	PSNR	LPIPS	#G	FPS	SSIM	PSNR	LPIPS	#G	FPS	SSIM	PSNR	LPIPS	#G	FPS	SSIM	PSNR	LPIPS	#G	FPS
Switch-NeRF	0.579	21.54	0.474	—	<0.1	0.654	22.57	0.457	—	<0.1	0.795	26.52	0.360	—	<0.1	0.541	23.62	0.609	—	<0.1
VastGaussian	<b>0.804</b>	<b>23.50</b>	—	—	—	<b>0.852</b>	<b>24.25</b>	—	—	—	<b>0.885</b>	<b>26.81</b>	—	—	—	<b>0.816</b>	<b>26.00</b>	—	—	—
CityGaussian (no LOD)	0.784	21.96	<b>0.243</b>	<b>13.30</b>	37.6	0.813	22.00	<b>0.211</b>	<b>10.80</b>	41.0	0.837	21.39	<b>0.230</b>	<b>3.80</b>	82.3	0.666	19.61	0.403	<b>16.41</b>	<b>35.1</b>
Hierarchical-3DGS (no LOD)	0.720	20.55	0.270	14.79	36.4	0.753	19.85	0.230	13.68	39.9	0.792	19.85	0.273	9.13	31.5	0.741	22.66	<b>0.297</b>	29.32	13.9
3DGS	0.787	22.42	0.282	<b>13.02</b>	<b>64.2</b>	0.807	21.96	0.256	<b>7.06</b>	<b>102.2</b>	0.833	21.26	0.285	<b>2.28</b>	<b>172.4</b>	0.718	19.83	0.370	<b>10.55</b>	18.5
Ours (no LOD)	<b>0.808</b>	<b>24.12</b>	<b>0.219</b>	18.22	<b>62.5</b>	<b>0.845</b>	<b>24.93</b>	<b>0.201</b>	14.15	<b>71.5</b>	<b>0.876</b>	<b>27.78</b>	<b>0.190</b>	8.39	<b>86.9</b>	<b>0.788</b>	<b>26.63</b>	<b>0.278</b>	42.42	<b>38.5</b>
CityGaussian	0.769	21.75	0.257	<b>3.49</b>	<b>83.6</b>	0.805	21.90	<b>0.217</b>	<b>3.13</b>	<b>65.7</b>	0.833	21.34	0.232	<b>1.77</b>	<b>113.4</b>	N/A(encountered a bug)				
Hierarchical-3DGS	0.695	20.18	0.296	6.59	46.3	0.741	19.70	0.243	10.01	44.2	0.788	19.82	0.278	6.67	36.0	<b>0.724</b>	<b>22.43</b>	<b>0.316</b>	<b>10.55</b>	<b>18.5</b>
Ours	<b>0.799</b>	<b>24.03</b>	<b>0.233</b>	<b>5.18</b>	<b>82.0</b>	<b>0.818</b>	<b>24.32</b>	<b>0.232</b>	<b>4.09</b>	<b>90.9</b>	<b>0.859</b>	<b>27.09</b>	<b>0.208</b>	<b>2.79</b>	<b>99.4</b>	<b>0.778</b>	<b>26.41</b>	<b>0.293</b>	<b>5.68</b>	<b>93.7</b>

Table 8. **Quantitative evaluation on Building, Residence, Sci-Art and Campus.** The results for VastGaussian are only partially available as it is not open-sourced and can only be obtained from its paper. All missing results are denoted by a “—”.

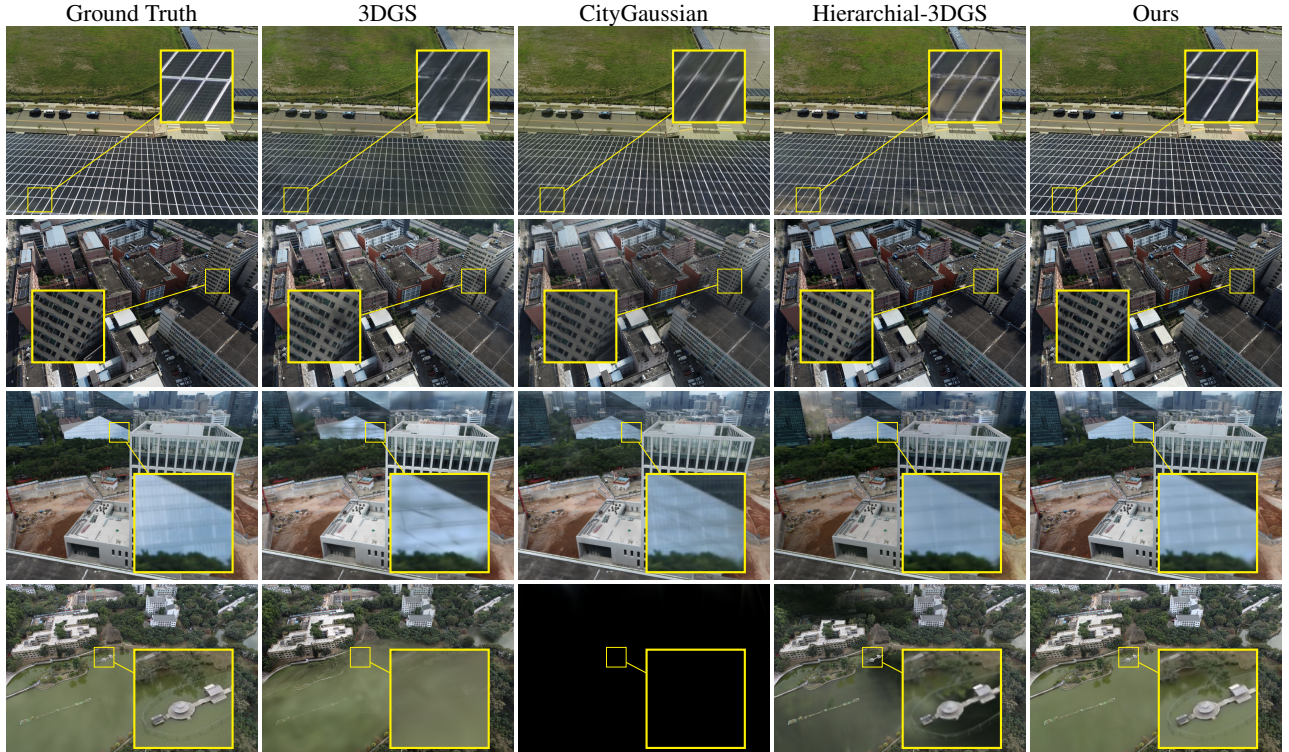


Figure 9. **Visualization results on Building, Residence, Sci-Art and Campus of ours and previous work.** All methods, except for 3DGS, render in LOD mode. The LOD mode of CityGaussian encountered a bug in the *Campus*, resulting in a completely black rendered image.

The second row of Figure 11 presents a statistical analysis of the similarity among Gaussians within a small local region, where 513 Gaussians are selected, and the similarities between 512 of them and a central reference Gaussian are computed to generate a histogram. In the absence of similarity regularization, most Gaussians exhibit low similarity, clustering around 0.1. Such low similarity results in significant differences in appearance transformations among Gaussians, leading to visible artifacts. In contrast, with similarity regularization applied, the similarity values among Gaussians predominantly exceed 0.8. This high degree of similarity ensures more consistent appearance adjustments across Gaussians, effectively preventing

the emergence of artifacts.

Scene	<i>Rubble</i>				<i>JNU-ZH</i>				<i>BigCity</i>			
Metrics	SSIM	PSNR	LPIPS	#G	SSIM	PSNR	LPIPS	#G	SSIM	PSNR	LPIPS	#G
1st level	0.870	28.34	0.153	3.71	0.889	26.57	0.114	5.41	0.925	27.48	0.098	10.38
2nd level	0.825	27.16	0.224	7.13	0.835	25.74	0.197	13.99	0.855	26.10	0.198	30.47
3rd level	0.826	27.29	0.228	13.52	0.822	25.85	0.232	25.58	0.847	26.62	0.219	75.15

Table 9. Quantitative evaluation of all the levels of our method, evaluated using the same downsampling factor as during training.

Scene	<i>Rubble</i>					<i>JNU-ZH</i>					<i>BigCity</i>				
Metrics	SSIM	PSNR	LPIPS	#G	FPS	SSIM	PSNR	LPIPS	#G	FPS	SSIM	PSNR	LPIPS	#G	FPS
w/o anti-aliasing	<b>0.817</b>	26.85	<b>0.237</b>	<b>3.02</b>	<u>103.5</u>	<b>0.817</b>	25.69	<b>0.233</b>	<b>4.92</b>	<u>67.8</u>	<b>0.847</b>	<b>26.52</b>	<b>0.213</b>	<b>6.68</b>	<u>73.6</u>
w/o absgrad	0.795	26.67	0.275	<u>3.15</u>	<b>109.8</b>	0.808	25.53	0.251	6.78	<b>68.3</b>	0.831	26.28	0.244	<u>6.74</u>	<b>78.1</b>
w/o tile-based cull.	<u>0.814</u>	<b>27.03</b>	<u>0.245</u>	3.60	83.0	<u>0.816</u>	<b>25.71</b>	<u>0.240</u>	<u>6.65</u>	54.6	<u>0.838</u>	<u>26.41</u>	<u>0.231</u>	6.84	65.8
full	<u>0.814</u>	<b>27.03</b>	<u>0.245</u>	3.60	99.7	<u>0.816</u>	<b>25.71</b>	<u>0.240</u>	<u>6.65</u>	63.9	<u>0.838</u>	<u>26.41</u>	<u>0.231</u>	6.84	73.0

Table 10. Additional qualitative ablations.

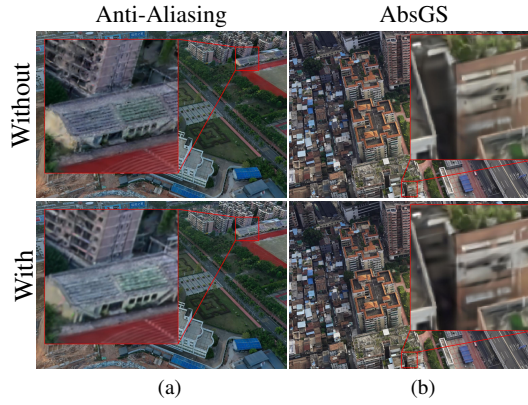


Figure 10. Visualization results of ablation on Anti-Aliasing and AbsGS.

Length	SSIM	PSNR	LPIPS
4	<u>0.816</u>	25.58	<u>0.237</u>
8	<u>0.815</u>	<u>25.62</u>	0.240
16	<b>0.822</b>	<b>25.85</b>	<b>0.232</b>

Table 11. The impact of the length of  $\ell^{(\mathcal{G})}$ .

Part. Size	SSIM	PSNR	LPIPS	FPS	#G ( $10^6$ )	#P
45m	0.803	25.27	0.256	35.9	<b>2.11</b>	860
90m	0.808	25.46	0.250	61.7	<u>3.06</u>	228
135m	<u>0.810</u>	<u>25.50</u>	<u>0.247</u>	<b>64.4</b>	3.68	140
180m	<b>0.813</b>	<b>25.62</b>	<b>0.243</b>	<u>63.1</u>	5.19	64

Table 12. Qualitative ablations of different rendering-time partition size on the *JNU-ZH* scene. #P represents the number of partitions.



Distances	SSIM	PSNR	LPIPS	FPS	#G ( $10^6$ )
(45m, 90m, $\infty$ )	0.789	25.00	0.271	<b>66.5</b>	<b>2.40</b>
(90m, 180m, $\infty$ )	0.808	25.46	0.250	<u>61.7</u>	<u>3.06</u>
(135m, 270m, $\infty$ )	<u>0.815</u>	<u>25.65</u>	<u>0.242</u>	58.8	3.62
(180m, 360m, $\infty$ )	<b>0.818</b>	<b>25.73</b>	<b>0.238</b>	56.2	4.13

Table 13. **Qualitative ablations of distance thresholds on the JNU-ZH scene with a partition size of 90m.** The distance values represent the maximum distances at which the 3rd, 2nd, and 1st LOD levels are used.

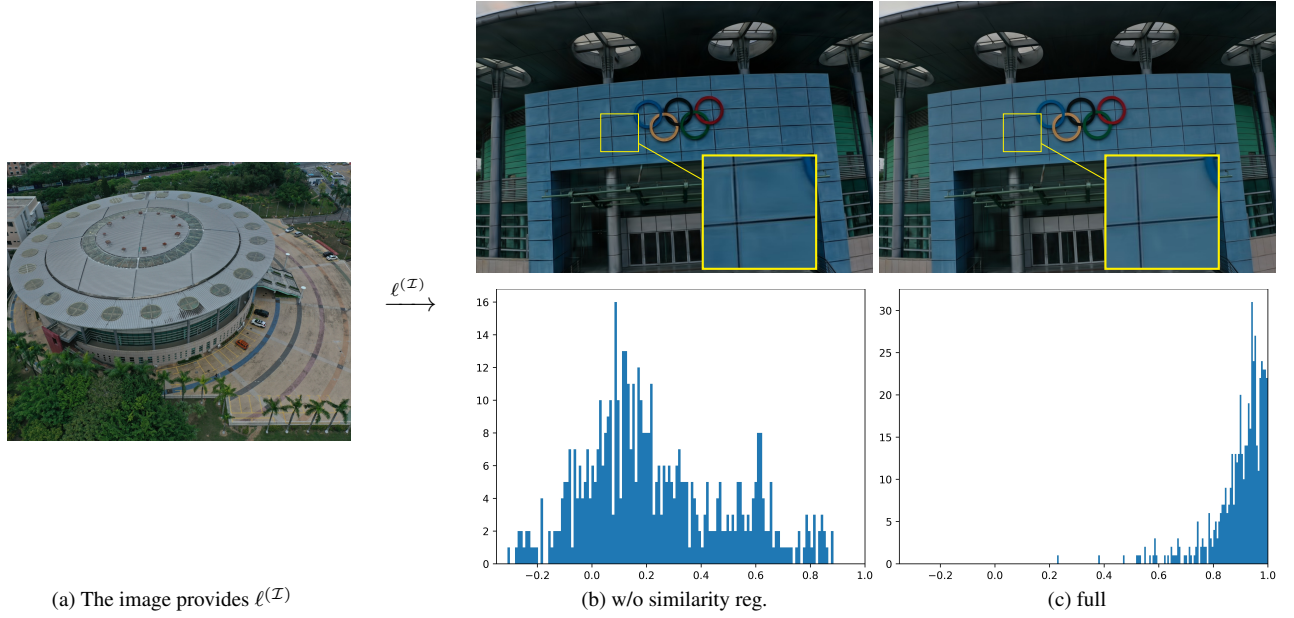


Figure 11. **A visual comparison of results with and without similarity regularization.**