

Radiant: Large-scale 3D Gaussian Rendering based on Hierarchical Framework

Haosong Peng^{*¶}, Tianyu Qi^{†¶}, Yufeng Zhan^{*||}, Hao Li[‡], Yalun Dai[§], Yuanqing Xia^{*}

^{*}School of Automation, Beijing Institute of Technology, Beijing, China

[†]School of Cyber Science and Technology, Sun Yat-sen University, Shenzhen, China

[‡]Brain and Artificial Intelligence Lab, Northwestern Polytechnical University, Xi'an, China

[§]Nanyang Technological University, Singapore

livion@bit.edu.cn, qity9@mail2.sysu.edu.cn, yu-feng.zhan@bit.edu.cn,

lifugan_10027@outlook.com, daiy0018@e.ntu.edu.sg, xia_yuanqing@bit.edu.cn

Abstract—With the advancement of computer vision, the recently emerged 3D Gaussian Splatting (3DGS) has increasingly become a popular scene reconstruction algorithm due to its outstanding performance. Distributed 3DGS can efficiently utilize edge devices to directly train on the collected images, thereby offloading computational demands and enhancing efficiency. However, traditional distributed frameworks often overlook computational and communication challenges in real-world environments, hindering large-scale deployment and potentially posing privacy risks. In this paper, we propose Radiant, a hierarchical 3DGS algorithm designed for large-scale scene reconstruction that considers system heterogeneity, enhancing the model performance and training efficiency. Via extensive empirical study, we find that it is crucial to partition the regions for each edge appropriately and allocate varying camera positions to each device for image collection and training. The core of Radiant is partitioning regions based on heterogeneous environment information and allocating workloads to each device accordingly. Furthermore, we provide a 3DGS model aggregation algorithm that enhances the quality and ensures the continuity of models' boundaries. Finally, we develop a testbed, and experiments demonstrate that Radiant improved reconstruction quality by up to 25.7% and reduced up to 79.6% end-to-end latency.

I. INTRODUCTION

3D Gaussian Splatting (3DGS) represents the latest advancement in the field of scene reconstruction of computer vision [1], [2]. Its ability to reconstruct real-world 3D structures directly from 2D images has led to widespread applications in virtual reality (AR/VR) [3], search-and-rescue [4], and autonomous driving systems [5], [6], all aiming to achieve superior performance. However, 3DGS-based methods pose challenges as they place a heavy load on the server and require extensive data storage and long completion time, especially when the scene is growing on a scale beyond the city scale.

Currently, many distributed 3DGS algorithms [7]–[9] have been developed for large-scale scenarios. A straightforward approach is cloud-based multi-GPU distributed training [7], as shown in the left part of Fig. 1, which has shown excellent results but requires substantial resources. Take VastGaussian [7] as an example, it requires two hours of training on $8 \times$ Tesla V100 GPU and 2.6 GB storage for an area of only

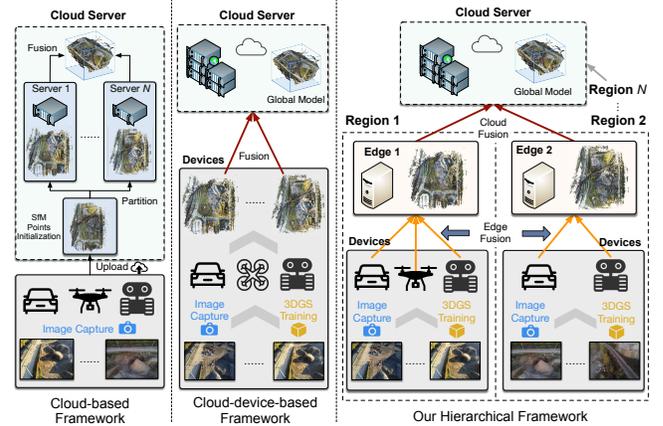


Fig. 1: **Left:** Cloud-based 3DGS training. **Middle:** Cloud-device-based 3DGS training. **Right:** Our hierarchical framework for 3DGS training.

$500 \times 250m^2$. Moreover, to obtain the initial model (i.e., the 3D Gaussian points) for the scene, all the raw images captured by the device must be uploaded to the cloud for structure-from-motion [10] (SfM), not only increasing communication overhead but also severely infringing on device privacy. Thanks to advancements in edge computing, cloud-device-based 3DGS has attracted attention [11], [12]. As shown in the middle part of Fig. 1, multiple devices simultaneously capture images in different regions. These images undergo local 3DGS training, and the models are then uploaded to the cloud for fusion, enabling comprehensive scene reconstruction tasks. However, cloud-device-based 3DGS algorithms do not address deployment issues in real-world scene reconstruction. The system's computational power and communication heterogeneity can cause stragglers. As the scale of the scene increases, the uploading of models by numerous devices can lead to network congestion, and the fusion of models in the cloud incurs significant overhead. Additionally, the cloud can access the precise location privacy of the devices [13]. Therefore, a more practical and efficient framework is urgently needed especially in large-scale scene reconstruction.

[¶]Equal Contribution.

^{||}Corresponding author: yu-feng.zhan@bit.edu.cn

Recently, the cloud-edge-device framework has been proposed to address deployment issues in large-scale systems [14]–[16]. When this framework is applied to distributed 3DGS, as shown in the right part of Fig. 1, each region can be effectively partitioned and trained. Specifically, the devices capture their images and independently train the models, which are then uploaded to the edge for preliminary fusion before being further fused in the cloud. In this way, the computational loads are ingeniously reallocated from the cloud to the edge, which can significantly mitigate the burden in the cloud. Besides, the images collected and the model trained by devices in different regions do not require sharing or direct uploading to the cloud, thus preserving data and location privacy. Moreover, this framework effectively enhances the scalability of the system. For new scenes, they can be developed as a new area by adding edges and devices to obtain a global model, without impacting the already trained regions.

Although combining the hierarchical framework with 3DGS can yield excellent results, the following challenges have to be addressed. In the real-world system, the pronounced heterogeneity can cause a severe straggler effect, ultimately affecting reconstructing efficiency [15]. For example, the varying number of images on each device leads to differences in model sizes, resulting in heterogeneous communication time for fusion. Even with the same number of images, training time differs across devices with varying computational power. On the other hand, the geographical locations of the edges and the camera positions are often constrained in real-world environments. Therefore, it is essential to consider how to partition regions geometrically for each edge and allocate workloads (i.e. camera positions) appropriately based on the heterogeneity of each device. Furthermore, we find that traditional merging methods cause obvious visual disruptions at the boundaries after model fusion. This impact worsens with the increase in the number of models merged. Developing an fusion scheme within the cloud-edge-device framework to maximize model effectiveness while addressing these issues remains a challenge.

In this work, we propose a hierarchical 3D Gaussian rendering framework for large scene reconstruction in a heterogeneous cloud-edge-device system, which jointly considers high efficiency, scalability, and privacy. We introduce a new 3DGS model aggregation algorithm within the hierarchical framework, which effectively addresses the shortcomings of traditional merging algorithms. On the other hand, we find that workload partitioning is crucial for the efficiency of modeling tasks in heterogeneous systems. Different from other distributed 3DGS algorithms, this algorithm fully considers the computational and communication heterogeneity within a three-layer architecture. Specifically, we propose an Adaptive Region Planning (ARP) algorithm based on geographical and resource characteristics to partition regions for edges, as well as a Resource-aware Task Partitioning (RTP) algorithm to allocate workloads for devices under each edge. We develop a prototype of Radiant and deploy it on a resource-heterogeneous testbed. Experiments on two real-world large-

scale scene datasets show that Radiant can achieve better reconstruction performance and lower latency compared to the state-of-the-art. Our contributions are as follows:

- To the best of our knowledge, we are the first to introduce the cloud-edge-device framework to large-scale scene reconstruction, and we design a model aggregation scheme tailored for this framework that maximizes model performance.
- We develop an optimization algorithm to allocate workloads based on the real-world distribution of edges and camera positions, as well as the system heterogeneity, to minimize the end-to-end latency.
- We develop a prototype of Radiant and design a testbed, with experiments indicating that Radiant can reduce end-to-end latency by up to 79.6% while achieving up to 25.7% better quality compared to the state-of-the-art.

The rest of the paper is organized as follows: Section II shows the background and motivation of this paper. Section III provides a detailed design of Radiant. We conduct extensive experiments in Section IV, and the related work is reviewed in Section V. Finally, we conclude this paper in Section VI.

II. BACKGROUND AND MOTIVATION

In this section, we first briefly review the background of 3D Gaussian Splatting. Then we conduct an experimental study to motivate the design of Radiant.

A. 3D Gaussian Splatting

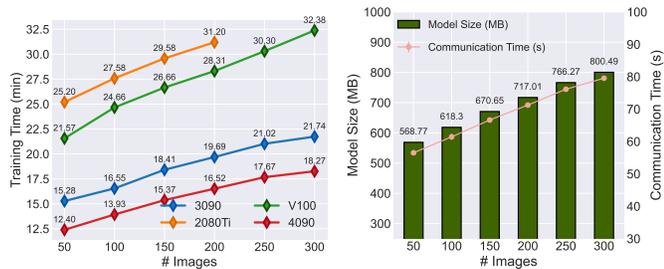
Our scene reconstruction is based on 3DGS [1] technology, where the geometry and appearance of the scene are rendered using a set of 3D Gaussian. The cameras that capture images can be written as $\mathbf{C} = \{(\mathcal{C}_j, I_j)\}$, where \mathcal{C}_j is the parameters of j -th camera including the position, intrinsic and extrinsic matrices, and I_j is the corresponding image. Let $\mathbf{G} = \{G_i\}$ be a set of 3D Gaussians. Each 3D Gaussian $G_i = (\mathbf{x}_i, \Sigma_i, \mathbf{f}_i, \alpha_i)$ is characterized by its position $\mathbf{x}_i = (x_i, y_i, z_i)$, covariance matrix $\Sigma_i \in \mathbb{R}^{3 \times 3}$, opacity $\alpha_i \in \mathbb{R}$, and spherical harmonic coefficients $\mathbf{f}_i \in \mathbb{R}^l$ for view-dependent colors, where l is related to the degree of the spherical harmonics. The initial position of 3D Gaussian is initialized using SfM [10]. 3DGS renders the novel view image $I_j^r = \mathcal{R}(\mathcal{C}_j, \mathbf{G})$ by a differentiable rasterizer \mathcal{R} .

Specifically, 3D Gaussians are projected into the 2D image space for rendering. The projected 2D covariance matrix Σ' is

$$\Sigma' = \mathbf{J}\mathbf{W}\Sigma\mathbf{W}^\top\mathbf{J}^\top, \quad (1)$$

where $\mathbf{W} \in \mathbb{R}^{3 \times 3}$ is the viewing transformation matrix and $\mathbf{J} \in \mathbb{R}^{3 \times 3}$ is the Jacobian of the affine approximation of the projective transformation [17]. Next, given the pixel \mathbf{p} and its position, we get a list of Gaussians \mathcal{N} sorted by their distance to the pixel. Then, the color C of each pixel is computed using α -blending, as

$$C = \sum_{i \in \mathcal{N}} c_i \alpha'_i \prod_{k=1}^{i-1} (1 - \alpha'_k), \quad (2)$$



(a) Training time v.s. # Images (b) Model size and comm. time.

Fig. 2: Heterogeneity when training across different scene scales on various devices.

where c_i is the learnable color computed from the spherical harmonics with f_i , and the final opacity α'_i is the multiplication result of the learnable opacity α_i and the Gaussian, as

$$\alpha'_i = \alpha_i \times \exp\left(-\frac{1}{2}(\mathbf{p}' - \mathbf{x}'_i)^\top \Sigma_i^{-1}(\mathbf{p}' - \mathbf{x}'_i)\right), \quad (3)$$

where \mathbf{p}' and \mathbf{x}' are coordinates in the projected space.

The properties of 3D Gaussians are optimized with respect to the loss function between I_j^r and the ground truth image I_j , as

$$\mathcal{L} = (1 - \lambda)\mathcal{L}_1(I_j^r, I_j) + \lambda\mathcal{L}_{D-SSIM}(I_j^r, I_j), \quad (4)$$

where λ is a hyper-parameter, and \mathcal{L}_{D-SSIM} denotes the D-SSIM loss. Finally, a set of trained 3D Gaussians \mathbf{G} is referred to as a model.

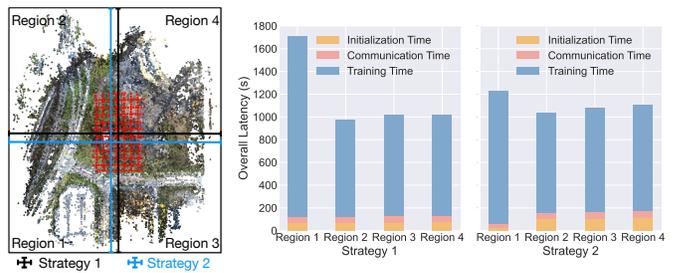
B. The Influence of Heterogeneity in 3DGS Training

To explore the heterogeneity of 3DGS local device training in the real world, we conduct multiple experiments to verify the impact of varying image quantities.

Device training time. To validate the influence on local model training time, we use four different devices with varying computational power: NVIDIA 4090, 3090, 2080Ti, and V100 to run the 3DGS training with varying quantities of images. Fig. 2a illustrates training time across various scene scales with different devices. Generally, as the number of images increases, the scene construction of the model becomes more complex, leading to longer overall training time. It can be observed that device training time significantly vary across devices with different computational powers.

Model size and communication time. We use the devices above to measure the model size and communication time under 5G bandwidth [18] after local training, as shown in Fig. 2b. Since more images result in more 3D Gaussian being generated, leading to larger model sizes, there are still significant communication time differences even under 5G conditions, causing some stragglers.

The results show that the varying quantities of images on devices significantly impact local training time. This variation also results in models of different sizes, which subsequently affects the communication time to the edge and cloud for fusion. Therefore, an effective workload partitioning algorithm



(a) Partitioning strategy (b) The overall latency of each region when equally partitioning and balanced partitioning.

Fig. 3: The overall latency breakdown of each device in two partitioning strategies.

is needed to ensure synchronous training across devices with various computational powers.

C. The Influence of Geographical Region Partitioning

In real-world scenarios, we need to assign nearby camera positions for image collection to each edge, thus forming a region. Devices within each region are designated to collect corresponding images. Due to the varying number and computational power of devices in each region associated with an edge, different regions result in heterogeneous completion times.

To validate the importance of region partitioning, we test the overall latency using different partitioning strategies in the same scenario. As shown in Fig. 3a, we deploy the low computational power devices in Region 1, while the high computational power devices were deployed in the other three regions. We establish two partitioning strategies for the four regions. Partitioning strategy 1 represents evenly distributing camera positions across four regions, while strategy 2 reduces the number of cameras in Region 1 and offloads them to other regions. The overall latency is illustrated in Fig. 2b, where it can be observed that Region 1 takes significantly longer. On the contrary, strategy 2 has a 28% latency reduction compared to the even distribution. Nevertheless, the impact of the two partitioning strategies on model quality (i.e., PSNR) is very minimal.

The above experiments motivate that a reasonable geographical partition of regions can mitigate heterogeneity among edges. However, designing an optimal and refined partitioning algorithm for complex systems remains a challenge that needs to be addressed.

D. The Influence of Different Fusion Scheme

In distributed 3DGS systems, the common model merging algorithm [7], [9] directly cuts the model boundary and concatenates them together. However, we observed that the merging method results in visual discontinuities at the boundaries (Fig. 12 visualizes this phenomenon). To conduct a quantitative study, we select a scene and divide it into multiple regions for training. We employ two fusion algorithms, a traditional merging method and a novel approach involving the upload of

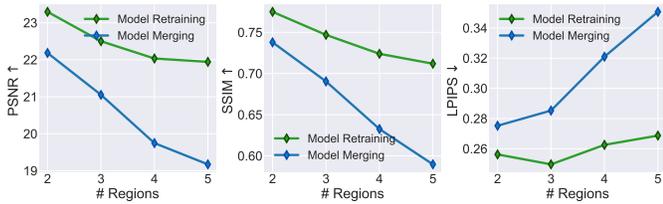


Fig. 4: The performance (\uparrow PSNR, \uparrow SSIM, \downarrow LPIPS) of using model merge and model retrain algorithms in the same area.

ground truth images followed by minimal retraining epochs of the entire model. As shown in Fig. 4, we evaluate three quality metrics of the global fused model. As the number of regions increases, the performance of the merging method deteriorates, while retraining method results in negligible performance loss.

The above results indicate that the merging method performs poorly due to boundary discontinuity, whereas retraining effectively mitigates this issue. However, retraining requires access to the ground truth images, which compromises the privacy of devices. Therefore, it is crucial to develop a fusion algorithm that utilizes retraining while preserving privacy.

III. RADIANT DESIGN

In this section, we first present the overview of the Radiant framework. Then, the workload partitioning problems and algorithms are introduced, followed by the model aggregation algorithm.

A. Overview

Radiant employs a cloud-edge-device hierarchical architecture to achieve efficient scene reconstruction. As illustrated in Fig. 5, consider a large-scale scene with a set of camera positions where images will be collected. This scene includes base stations at different positions, each serving as an edge center. Within each edge, there are multiple devices with various computational powers. All the edges and devices should collaborate to reconstruct the entire scene.

The workflow of Radiant is introduced as follows: (a) The large scene has to be divided into multiple regions, with its corresponding camera positions allocated to the edge. The edge then distributes these camera positions among various devices for image collection and training. (b) Then, all the devices independently collect their images and initialize the 3D Gaussian point by SfM, followed by training their local model. (c) After all devices in the same edge complete their training, the models are uploaded to the edge for aggregation. (d) All edge models are then uploaded to the cloud for the final aggregation, finally producing a global model of the scene.

Our design adheres to the following objectives: (1) **Efficiency**. We allocate the workload across each device appropriately to minimize completion time while enhancing performance in the large-scale system. (2) **Scalability**. Our framework allows new scenes to be added at any time, and once trained, these can be integrated with the models of previous scenes. (3) **Privacy**. Images captured by devices

are not uploaded to the edge or cloud, and there is no data exchange between devices. Additionally, the cloud does not have access to the location information of each device.

Specifically, to balance the workload on each edge, we propose the ARP algorithm to allocate the number and positions of cameras based on each region’s computational power and edge’s bandwidth (Sec. III-B). For device workload partitioning, we introduce the RTP Algorithm, which allocates the camera positions to each device based on computational power to mitigate straggler effects (Sec. III-C). Then, all the devices can conduct 3DGS training in parallel. Next, we introduce a model aggregation algorithm designed to achieve high quality while ensuring rapid execution and privacy. Finally, all the edge models are uploaded to the cloud and fused to form a global model (Sec. III-D).

B. Adaptive Region Planning Algorithm

Suppose there are M edges $\mathcal{M} = \{1, 2, \dots, M\}$ located at position $\mathcal{P} = \{P^1, \dots, P^M\}$. Edge m has a set of devices $\mathcal{N}^m = \{1, 2, \dots, N^m\}$. The communication bandwidth between each edge m and the cloud is denoted as B^m , and the device bandwidth is B^n . Therefore, the task completion time T^m for each edge m includes the training time T_{train}^m , the aggregation time T_{aggre}^m and the communication time between edge and cloud T_{com}^m , which is

$$T^m = T_{train}^m + T_{aggre}^m + T_{com}^m. \quad (5)$$

T_{com}^m is calculated as the edge model size Z^m divided by the average bandwidth, defined as

$$T_{com}^m = \frac{Z^m}{B^m}, \quad (6)$$

where Z^m can be approximated as the sum of the size of device models, i.e., $Z^m \approx \sum_{n=1}^{N^m} Z^n$. The aggregation time T_{aggre}^m is the time taken for all device models to be fused on the edge. In our method, this is essentially equivalent to model retraining, therefore, it is proportional to the number of images captured under each region. The aggregation time of edge m is calculated as:

$$T_{aggre}^m = \mathcal{F}_m(|\mathbf{C}^m|), \quad (7)$$

where $|\mathbf{C}^m|$ is the amount of camera positions that assigned to edge m and \mathcal{F}_m is a function that can be fitted. Within each edge, all devices can start training simultaneously, so the total completion time is determined by the device that takes the longest to complete its training:

$$T_{train}^m = \max_{n=1}^{N^m} \{T^n\}. \quad (8)$$

The task completion time T^n of device n is related to the number of camera positions allocated to it and its computational power, which can be calculated using the RTP algorithm introduced in Sec. III-C. In summary, to minimize the total completion time of the modeling task, we need to minimize the maximum completion time of any edge. This optimization problem can be formulated as

$$\min \max_{m=1}^M \{T_{train}^m + T_{aggre}^m + T_{com}^m\}, \quad (9)$$

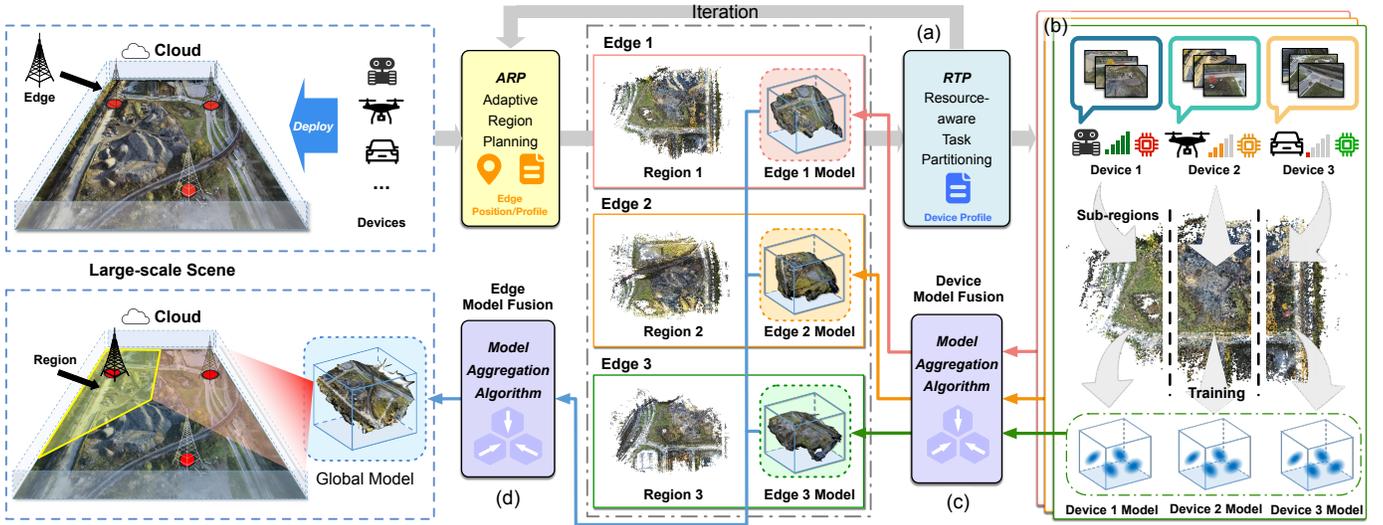


Fig. 5: The Overview of Radiant. (a) The Adaptive Region Planning algorithm iterates together with the Resource-aware Task Partitioning algorithm. (b) Device training. (c) Device model aggregation. (d) Edge model aggregation.

$$\text{s.t. } \bigcup_{m=1}^M \mathbf{C}^m = \mathbf{C}, \quad (10)$$

$$\mathbf{C}^i \cap \mathbf{C}^j = \emptyset, \quad \forall i, j \in \{1, 2, \dots, M\}, i \neq j, \quad (11)$$

(6), (7), (8).

Eqn. (10) indicates that the region of each edge needs to cover all camera positions, while Eqn. (11) specifies that the regions must not overlap, thereby preserving the privacy of each edge.

The basic idea of the ARP is to achieve load balancing across all edges, which ensures that all edges complete their tasks in approximately the same amount of time. Suppose we have the coordinates of all edge centers, the computational power of each device, and the positions of all images that need to be collected. The output is the final region of each edge and its camera position list. The pseudo-code is shown in Algorithm 1, which contains the following steps.

- 1) **Initializing the regions.** We generate the initial division of N regions (Line 1) following the rule of the Voronoi diagram [19]. In the Voronoi diagram, the distance from all camera positions within each region to their respective edge center \mathcal{P} is shorter than any other edge center. This partitioning ensures that devices can collect images in the nearest area, enhancing efficiency in the real systems.
- 2) **Estimate the completion time.** The camera positions within each region are incorporated into the corresponding edge (Line 2). We estimate the completion time for all edges using Eqn. (5), then calculate their average time \bar{T} (Lines 3-4).
- 3) **Move the boundary of the outlier.** We select the edge with the greatest deviation from the average estimated time, determined by:

$$m^* = \arg \max_m \{|T^m - \bar{T}|\}. \quad (12)$$

Algorithm 1: Adaptive Region Planning Algorithm

Input: The positions of edges $\mathcal{P} = \{P^1, \dots, P^M\}$, Camera position \mathbf{C} , distance d ;

Output: Regions $\mathbf{R} = \{r^1, \dots, r^M\}$; Camera position lists $\mathbf{L} = \{l^1, \dots, l^M\}$

```

1  $\mathbf{R} \leftarrow \text{Voronoi}(\mathcal{P})$  # Initialize the regions;
2 for  $m = 1, 2, \dots, M$  do
3    $\mathbf{C}^m \leftarrow \text{Find}(\mathbf{C}, r^m)$  # collect camera positions;
4    $T^m \leftarrow T_{train}^m + T_{aggre}^m + T_{com}^m$ ;
5 end
6  $\bar{T} \leftarrow \frac{1}{M} \sum_{m=1}^M T^m$ ;
7 while  $\max_m \{|T^m - \bar{T}|\} > T_{thres}$  do
8    $m^* = \arg \max_m \{|T^m - \bar{T}|\}$ ;
9    $\mathbf{R} \leftarrow \text{Move\_boundary}(m^*, \mathbf{R}, d)$ ;
10  for  $m = 1, 2, \dots, M$  do
11     $\mathbf{C}^m \leftarrow \text{Find}(\mathbf{C}, r^m)$ ;
12     $T^m \leftarrow T_{train}^m + T_{aggre}^m + T_{com}^m$ ;
13  end
14   $\bar{T} \leftarrow \frac{1}{M} \sum_{m=1}^M T^m$ ;
15 end
16 for  $m = 1, 2, \dots, M$  do
17    $l^m.append(\mathbf{C}^m)$ ;
18 end

```

If T^{m^*} of edge m^* exceeds the average \bar{T} , the boundary is shrunk by a distance d in the direction of the boundary's normal; if it is below the average, the boundary is expanded by d in the normal direction (Line 7).

- 4) **Repeat until the condition is satisfied.** At each round, the completion time of each edge for the new regions can be re-estimated, and the process is repeated until the residual completion time of every edge is below a

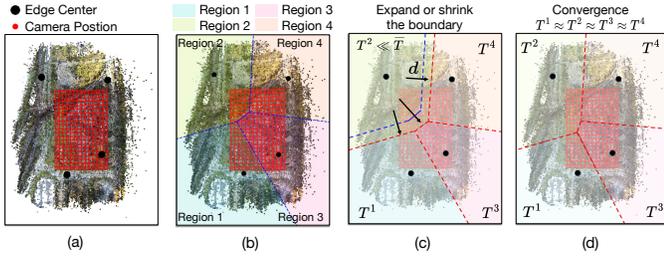


Fig. 6: An example of Adaptive Region Planning algorithm. (a) Original map of camera positions and edge centers. (b) Initialize the regions with a Voronoi diagram. (c) Iteratively adjusts the boundaries of the region to balance the load. (d) Converges to equilibrium.

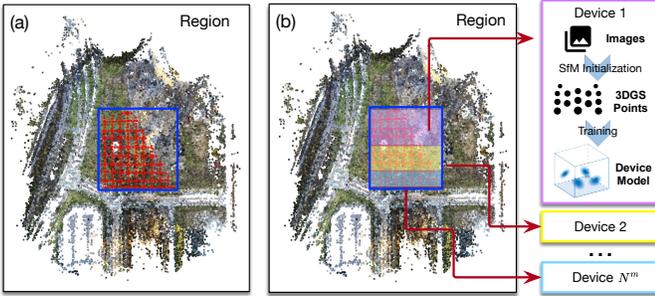


Fig. 7: An example of Resource-aware Task Partitioning algorithm.

specified threshold (Lines 5-11).

Fig. 6 illustrates an example of the ARP algorithm. The four edge centers, marked by black circles, are situated at distinct positions, while the red dots represent the camera position. We initially partitioned these into four regions using a Voronoi diagram. However, this allocation results in a highly imbalanced workload on each edge because of the heterogeneity of devices. For example, region 2 is allocated very few cameras and consequently has a very fast completion time, which is identified as an outlier. Therefore, we expand its boundary by d units and the completion time can be re-estimated. We continue to iterate this process until the completion times of all four edges converge towards the same. Finally, we record the positions and amount of the cameras assigned to each edge and the boundaries of their respective regions. The ARP algorithm benefits from favorable initial conditions and can quickly converge through rapid iterations, efficiently allocating regions to edges without overlap within several seconds.

C. Resource-aware Task Partitioning Algorithm

At each iteration of the ARP algorithm, we need to estimate the optimal completion time for each edge to complete its tasks. This problem can be modeled as a load-balancing assignment issue, with the objective of minimizing the longest completion time among all devices. We already know the partitions of camera position quantities to each edge and need to allocate them to the resource-heterogeneous devices next.

For simplicity, we only consider the situation in edge m as the problem is identical within each edge. Specifically, the completion time for each device n equals the sum of the initialization time, training time, and the communication time between the device and the edge, which is $T^n = T_{init}^n + T_{train}^n + T_{com}^n$. The initialization time for 3D Gaussian is related to the image number used by the device, which can be profiled by

$$T_{init}^n = \mathcal{G}(|\mathbf{C}^n|). \quad (13)$$

The communication time is calculated as the device model size Z^n divided by the average bandwidth, and the model size is related to the number of cameras, defined as:

$$Z^n = \mathcal{H}(|\mathbf{C}^n|), \quad (14)$$

where function $\mathcal{H}(\cdot)$ can be profiled offline, as illustrated in Fig. 2b. Therefore, T_{com}^n is calculated by

$$T_{com}^n = \frac{Z^n}{B^n}. \quad (15)$$

Furthermore, we can determine the function between the training time and different cameras for devices with various computational powers by offline profiling as illustrated in Fig. 2a. Hence, the training time is

$$T_{train}^n = \mathcal{F}_n(|\mathbf{C}^n|), \quad (16)$$

Overall, this problem can be formulated as:

$$\min \max_{n=1}^{N^m} \{T_{init}^n + T_{train}^n + T_{com}^n\}, \quad (17)$$

$$\text{s.t.} \quad \bigcup_{n=1}^{N^m} \mathbf{C}^n = \mathbf{C}^m, \quad (18)$$

$$\mathbf{C}^i \cap \mathbf{C}^j = \emptyset, \quad \forall i, j \in \{1, 2, \dots, N\}, i \neq j, \quad (19)$$

$$|\mathbf{C}^n| \leq |\mathbf{C}_{max}^n|. \quad (20)$$

$$(13), (15), (16)$$

Eqn. (20) stipulates that the number of camera positions allocated to each device must not exceed the maximum capacity that its video memory can handle. The optimal allocation will be achieved when the completion times of all devices are equal. Consider that \mathcal{G} , \mathcal{F} , and \mathcal{H} are known, the number of camera positions allocated to each device can be directly calculated. The result serves as the estimated values for T_{train}^m from Eqn. (8), which contributes to the iterations of the ARP algorithm.

Fig. 7 presents an example of the RTP algorithm. Without loss of generality, we divide the camera positions into N^m sub-regions along the longer axis of the rectangular area they occupy according to any device order. The number of images captured by cameras for each sub-region is determined by the solution to problem (17). Finally, each device begins by initializing the Gaussian points using the captured images, followed by 3DGS training for reconstructing a specific sub-region. The combined scope of these sub-regions is sufficient to cover the entire region of the edge.

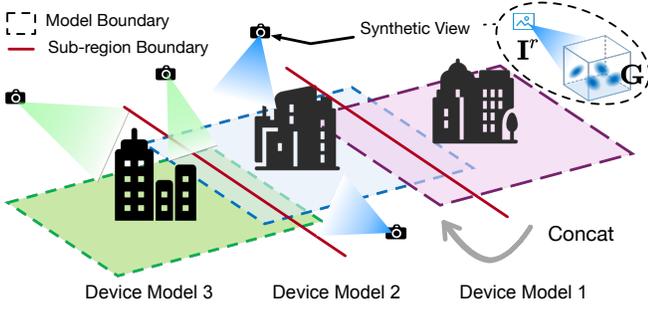


Fig. 8: Model Aggregation. We concatenate device models directly to form an edge model and retrain it using synthetic views.

Through the ARP and RTP algorithms, the reconstruction tasks for the entire scene are meticulously distributed to each device. This ensures that all devices can complete their respective tasks as efficiently as possible while mitigating the effects of stragglers.

D. Model Aggregation Algorithm

Traditional cutting and merging according to region boundaries does not yield an ideal reconstruction quality. To address this issue, we propose a model aggregation algorithm for the edge that seamlessly stitches together all device models. Specifically, each device model is trained using non-overlapping cameras C^n . Synthetic views I^r can be rendered based on the perspectives of camera C^n after the devices upload their model to the edge, which can be expressed as $I^r = \mathcal{R}(C^n, G^n)$, where G^n is the model of device n . We concatenate all device models together and retrain the edge model for a few epochs using these synthetic views I^r as depicted in Fig. 8. This approach can significantly enhance the quality of the model boundaries while eliminating the need to upload images to the edge, protecting the privacy of devices.

The aggregation of cloud models follows the same steps. In cases where cloud resources are insufficient for aggregation, we can still employ the merging method. Specifically, according to the final region boundaries recorded in the ARP algorithm, we cut out the 3D Gaussians outside the boundary lines of each region from the edge models, and then directly concatenate them.

IV. EVALUATION

We first describe the implementation of Radiant and experimental settings, followed by Radiant performance and its latency comparison. Finally, we conduct ablation studies and visualization to prove the effectiveness of the Radiant.

A. Implementation

Our Radiant framework is implemented with Python. We use the HUAWEI Elastic Cloud Server (ECS) [20] with 8 vCPU and 16G memory as the cloud server. In addition, we use the following servers with different GPUs as the heterogeneous device set. **Instance #1**: A desktop equipped

TABLE I: Composition of the heterogeneous systems.

	System 1	System 2
Cloud	HUAWEI Cloud ECS	HUAWEI Cloud ECS
Edge	1 × 4090	1 × 4090
Region 1	Instance #1	Instance #1
Region 2	Instance #1 & Instance #4	Instance #1 & Instance #4
Region 3	Instance #2 & Instance #3	Instance #2 & Instance #3
Region 4	-	Instance #2 & Instance #4

with two Intel(R) Xeon(R) Silver 4310 CPUs, 128 GB of RAM, and two NVIDIA 4090 GPUs, each with 24 GB of VRAM. **Instance #2**: A desktop equipped with an Intel(R) Xeon(R) Silver 4310 CPU, 377 GB of RAM, and eight NVIDIA V100 GPUs, each with 32 GB of VRAM. **Instance #3**: A desktop equipped with an Intel(R) Core(TM) i5-9400F CPU, 24 GB of RAM, and an NVIDIA 2080Ti GPU with 11 GB of VRAM. **Instance #4**: A desktop equipped with an Intel(R) Core(TM) i7-13700KF CPU, 64 GB of RAM, and an NVIDIA 3090Ti GPU with 24 GB of VRAM. We use 5G bandwidth data in our system [18].

B. Experiment Settings

1) *Performance Metrics*: The performance metrics we evaluate include end-to-end latency and global model quality. We assess the quality of scene reconstruction using three quantitative metrics: PSNR, SSIM [21] and VGG-based LPIPS [22].

2) *Datasets*: We use the Mill 19 dataset [23] which consists of *Rubble* and *Building* scenes. The *Rubble* and *Building* datasets consist of 1,678 and 1,940 images respectively with a resolution of 4608×3456 . The two scenes together cover a large-scale area of approximately $100,000 m^2$. We resize them to 1152×864 for training and evaluation, following the same setting of [7].

3) *Training Hyper-parameters*: Each device is optimized for 20,000 iterations. The densification [24] starts at the 500-th iteration and ends at the 15,000-th iteration, with an interval of 100 iterations. The number of epochs for retraining during model aggregation is set to 10. The other settings are identical to those of 3DGS [24].

4) *Baselines*: We compare Radiant with centralized and distributed training methods. We will also include the NeRF-based [25] 3D reconstruction method in our comparison. The centralized approach includes GP-NeRF [26], and 3DGS [24], Switch-NeRF [27], and the distributed method includes Vast-Faussian [7], Mega-NeRF [23], Drone-NeRF [28], Fed-NeRF [29], and Fed3DGS [11]. Additionally, we present the privacy and scalability features of each method to ensure a fair comparison.

C. Radiant Performance

To demonstrate the performance of Radiant in large heterogeneous systems, we set up two groups of heterogeneous systems to complete scene reconstruction tasks in the *Rubble* and *Building* scene, respectively. Table I shows the composition of the two systems we used, the number of edges,

TABLE II: Performance of Radiant and several baselines.

Architecture	Method	Rubble			Building			Scalability	Privacy	
		PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow			
Centralized Framework	GP-NeRF [26]	24.08	0.563	0.497	20.99	0.565	0.490	✗	✗	
	3DGS \diamond [24]	25.51	0.725	0.316	22.53	0.738	0.214	✗	✗	
	Switch-NeRF [27]	24.31	0.562	0.496	21.54	0.579	0.474	✗	✗	
Distributed Framework	Cloud-based	Drone-NeRF [28]	19.51	0.528	0.489	18.46	0.490	0.469	✓	✗
		Mega-NeRF [23]	24.06	0.553	0.516	20.93	0.547	0.504	✓	✗
		VastGaussian [7]	25.20	0.742	0.264	21.80	0.728	0.225	✗	✗
	Cloud-device-based	FedNeRF [29]	20.12	-	-	17.51	-	-	✓	✗
		Fed3DGS [11]	20.62	0.588	0.437	18.66	0.602	0.362	✓	✗
	Hierarchical	Even-system1 \triangle	23.93	0.719	0.324	21.19	0.682	0.304	✓	✓
		Even-system2 \triangle	24.08	0.729	0.310	21.33	0.698	0.297	✓	✓
		Radiant-system1	24.44	0.737	0.305	21.57	0.716	0.292	✓	✓
		Radiant-system2	24.53	0.740	0.284	21.66	0.723	0.236	✓	✓

\diamond modified 3DGS [24] so that it can be optimized on a single instance; \triangle “even” means to divide the images equally for each device.

and the configuration of edges. It is worth noting that each edge comprises several devices simulated by the **Instance** as mentioned in Sec IV-A. The locations of edge servers are randomly and uniformly distributed throughout the scene.

Table II presents the performance of Radiant alongside several baseline methods. We categorize the methods into centralized and distributed frameworks, with the latter further divided into cloud-based, cloud-device-based, and hierarchical frameworks. We also compare another Radiant version that evenly divides the workload to each device. As a novel method based on the hierarchical framework, Radiant outperforms all cloud-device framework methods in terms of PSNR, SSIM, and LPIPS. Compared to centralized and cloud-based architecture methods, it achieves comparable performance. However, the centralized method based on 3DGS requires collecting all images and initializing Gaussian points before training, thus failing to protect privacy. Similarly, it lacks scalability when new areas need to be planned and reconstructed. In FedNeRF and Fed3DGS, devices directly upload their models to the cloud, revealing each device’s location and thus failing to satisfy privacy requirements. Our Radiant method achieves high-quality reconstruction whether the workload is evenly distributed or divided according to heterogeneous systems. Specifically, Radiant can improve reconstruction quality (PSNR) by up to 25.7%. Besides, it satisfies both privacy and scalability requirements.

D. Latency Comparison

We compare the latency of Radiant with Fed3DGS [11] in Fig. 9. We implement the experiments using the Fe3DGS under the same scale of System 1 and System 2. Methods based on the NeRF usually require more than 1 day [7], [23] training, making them impractical for timely applications, and thus are excluded from our comparison. Specifically, Radiant is 19.27% and 20.99% faster than even distribution in the *Rubble* and *Building* scenarios, respectively, because even distribution slows down the completion time of the devices with the least

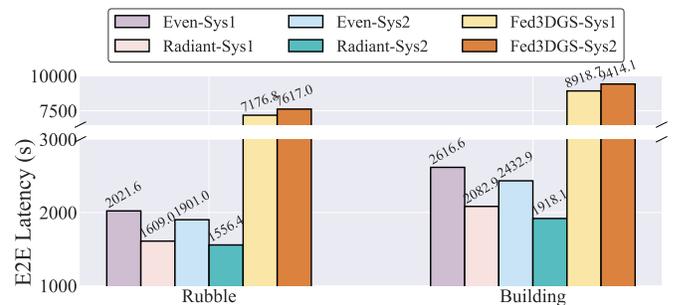


Fig. 9: Comparison of end-to-end latency. “even” means to divide the images equally for each device.

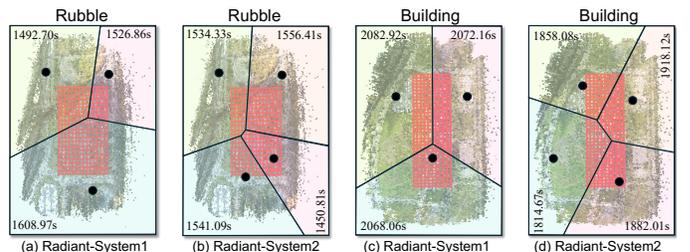


Fig. 10: The partitions of Radiant and the latency of each region.

computational resources. All our Radiant versions are more than 70% faster than Fed3DGS. For instance, our Radiant-system2 is 79.6% faster than Fededgs on *Building* dataset. This is due to Fed3DGS fusing all device modes sequentially on a single machine, which is highly time-consuming.

Furthermore, we demonstrated the partitioning results of the ARP algorithm under the settings of System 1 and System 2 in Fig. 10, as well as the latency on each edge. This figure illustrates that the ARP algorithm can effectively balance the workload distribution based on the computational resources of each region.

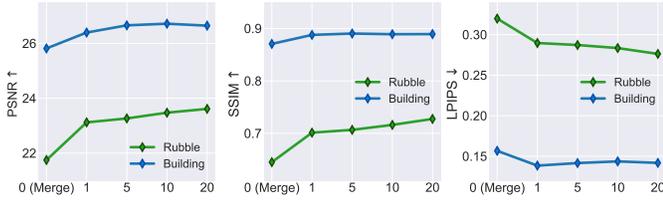


Fig. 11: The impact of the number of aggregate epochs on the quality of regional boundaries.



Fig. 12: The visualizations of two model fusion methods. The merging approach results in noticeable discontinuities at the boundaries.

E. Deep Dive into Radiant

We conducted ablation studies on the model’s aggregation methods. Table III shows a performance comparison between the model aggregation algorithm and directly merging according to the region’s boundaries. It is evident that the aggregated model maintains better accuracy across three metrics. In terms of the selection of retaining epochs, Fig. 11 demonstrates the impact of different training epochs on the quality of regional boundaries. We find that using only a small number of epochs can already achieve much better results than the merging method.

The differences between these merging and aggregation methods are more visually apparent, as shown in Fig. 12. We displayed the effects of two fusion methods with images from two different datasets. It was observed that using the merging method results in noticeable boundary lines, whereas the aggregation algorithm significantly alleviates this phenomenon.

We also display the convergence process of the ARP algorithm in Fig. 13. The y-axis represents the ratio of the maximum time difference in the current partitioned region to the average time, i.e., $(\max_m |T^m - \bar{T}|) / \bar{T}$. A smaller value indicates a more balanced workload distribution. Our ARP algorithm can rapidly converge after a small number of iterations.

V. RELATED WORK

In this section, we briefly review the latest work on 3DGS technology and large-scale hierarchical frameworks developed

TABLE III: Fusion methods ablation experiment.

Settings	Method	Rubble			Building		
		PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓
System1	Merge	22.51	0.637	0.371	21.17	0.709	0.278
		24.44	0.737	0.305	21.57	0.716	0.292
System2	Merge	23.00	0.682	0.328	20.46	0.658	0.310
		24.53	0.740	0.284	21.66	0.723	0.236

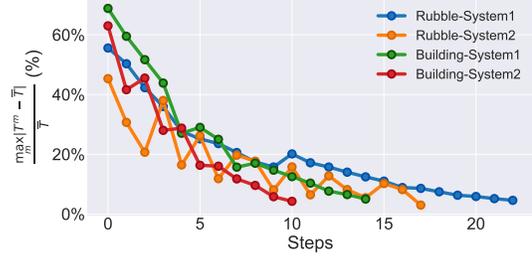


Fig. 13: The convergence process of the ARP algorithm.

in recent years.

A. 3D Scene Reconstruction

Traditional approaches [30]–[32] follow a structure-from-motion pipeline that estimates camera poses and generates sparse point clouds. However, such methods often contain artifacts or holes in areas with limited texture or specular reflections as they are challenging to triangulate across images. Recently, NeRF [25] and 3DGS [1] variants have become a worldwide 3D representation system thanks to their photo-realistic characteristics and the ability of novel-view synthesis, which inspires many works [7], [11], [23], [26], [27], [33], [34] to extend it into large-scale scene reconstructions. The above methods can be categorized into centralized and distributed frameworks. Centralized methods [26], [34] adopt the integration of NeRF-based and grid-based methods to model city-scale reconstruction. Distributed methods [7], [23] apply scene decomposition for multiple NeRF / Gaussian models optimization. However, with the growing scene size, both centralized and distributed variants limit their scalability due to the central server’s limited data storage and unacceptable computation costs. Nearly proposed Fed3DGS [11] introduces federated learning to leverage the computational resources of clients and emerge clients’ 3DGS models into the central server by a distillation update scheme. Nevertheless, this method ignores the scales and the bundle adjustment between different edge models, resulting in a performance drop. Meanwhile, it only focuses on over-fitting the photo-realistic rendering but ignores the geometry performance.

B. Large-Scale Hierarchical Framework

The hierarchical cloud-edge-device framework [14], [16], [35] has been proposed to address large-scale distributed edge computing problems and has been widely applied in various distributed domains in recent years. Wang et al. [15] identified the optimal clustering configuration and implemented

hybrid federated learning using a blend of synchronous and asynchronous methods. Li et al. introduced FedGS [36] for the IIoT environment, which utilizes the gradient-based binary permutation algorithm (GBP-CS) to select subsets. Mhaisen et al. [37] developed an optimization for the pairing of devices and edges to reduce the distance in class distribution. Zhong et al. [38] proposed the FLEE algorithm, which leverages edges and devices to achieve dynamic model updates. Deng et al. [39] proposed FedHKT, a hierarchical framework that utilizes a hybrid knowledge transfer mechanism to enhance learning efficiency and performance. Yang et al. proposed HierMo [40], which applies momentum to accelerate distributed training. Chen et al. [41] introduced SD-GT, which enhances federated learning on fog networks by optimizing device updates and improving model performance.

VI. CONCLUSION

We designed Radiant, a hierarchical framework for 3D Gaussian rendering in large scene reconstruction with resource-heterogeneous systems. It is the first distributed framework that simultaneously satisfies efficiency, scalability, and privacy. We propose the ARP and RTP algorithms to distribute the reconstruction workload among edges and devices. Additionally, we develop a novel model aggregation method for the hierarchical framework to enhance the model quality. Finally, we develop a testbed and experiments demonstrate that Radiant outperforms the state-of-the-art methods on model quality and end-to-end latency.

REFERENCES

- [1] 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.
- [2] T. Wu, Y.-J. Yuan, L.-X. Zhang, J. Yang, Y.-P. Cao, L.-Q. Yan, and L. Gao, “Recent advances in 3d gaussian splatting,” *Computational Visual Media*, pp. 1–30, 2024.
- [3] H. Li, J. Li, D. Zhang, C. Wu, J. Shi, C. Zhao, H. Feng, E. Ding, J. Wang, and J. Han, “Vdg: Vision-only dynamic gaussian for driving simulation,” *arXiv preprint arXiv:2406.18198*, 2024.
- [4] T. Zhang, K. Huang, W. Zhi, and M. Johnson-Roberson, “Darkgs: Learning neural illumination and 3d gaussians relighting for robotic exploration in the dark,” *arXiv preprint arXiv:2403.10814*, 2024.
- [5] P. Sun, H. Kretzschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine, et al., “Scalability in perception for autonomous driving: Waymo open dataset,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2446–2454, 2020.
- [6] H. Li, Y. Gao, D. Zhang, C. Wu, Y. Dai, C. Zhao, H. Feng, E. Ding, J. Wang, and J. Han, “Ggrt: Towards generalizable 3d gaussians without pose priors in real-time,” in *Proceedings of the European conference on computer vision (ECCV)*, 2024. Accepted.
- [7] 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*, pp. 5166–5175, 2024.
- [8] Y. Liu, H. Guan, C. Luo, L. Fan, J. Peng, and Z. Zhang, “Citygaussian: Real-time high-quality large-scale scene rendering with gaussians,” in *Proceedings of the European conference on computer vision (ECCV)*, 2024. Accepted.
- [9] G. H. L. Yu Chen, “Dogaussian: Distributed-oriented gaussian splatting for large-scale 3d reconstruction via gaussian consensus,” in *arXiv*, 2024.
- [10] N. Snavely, S. M. Seitz, and R. Szeliski, “Photo tourism: exploring photo collections in 3d,” in *ACM siggraph 2006 papers*, pp. 835–846, 2006.
- [11] T. Suzuki, “Fed3dgs: Scalable 3d gaussian splatting with federated learning,” *arXiv preprint arXiv:2403.11460*, 2024.
- [12] 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.
- [13] Y. Zhao and J. Chen, “Vector-indistinguishability: location dependency based privacy protection for successive location data,” *IEEE Transactions on Computers*, vol. 73, no. 4, pp. 970–979, 2023.
- [14] L. Liu, J. Zhang, S. Song, and K. B. Letaief, “Client-edge-cloud hierarchical federated learning,” in *ICC 2020-2020 IEEE international conference on communications (ICC)*, pp. 1–6, IEEE, 2020.
- [15] Z. Wang, H. Xu, J. Liu, H. Huang, C. Qiao, and Y. Zhao, “Resource-efficient federated learning with hierarchical aggregation in edge computing,” in *IEEE INFOCOM 2021-IEEE conference on computer communications*, pp. 1–10, IEEE, 2021.
- [16] Y. Cui, K. Cao, J. Zhou, and T. Wei, “Optimizing training efficiency and cost of hierarchical federated learning in heterogeneous mobile-edge cloud computing,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 42, no. 5, pp. 1518–1531, 2022.
- [17] M. Zwicker, H. Pfister, J. Van Baar, and M. Gross, “Ewa splatting,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 8, no. 3, pp. 223–238, 2002.
- [18] M. Ghoshal, Z. J. Kong, Q. Xu, Z. Lu, S. Aggarwal, I. Khan, Y. Li, Y. C. Hu, and D. Koutsonikolas, “An in-depth study of uplink performance of 5g mmwave networks,” in *Proceedings of the ACM SIGCOMM Workshop on 5G and Beyond Network Measurements, Modeling, and Use Cases*, pp. 29–35, 2022.
- [19] F. Aurenhammer, “Voronoi diagrams—a survey of a fundamental geometric data structure,” *ACM Computing Surveys (CSUR)*, vol. 23, no. 3, pp. 345–405, 1991.
- [20] HUAWEI CLOUD, “Huawei cloud elastic cloud server (ecs),” 2024. <https://www.huaweicloud.com/intl/en-us/product/ecs.html>, Last accessed on 2024-7-16.
- [21] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, “Image quality assessment: from error visibility to structural similarity,” *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [22] 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*, pp. 586–595, 2018.
- [23] 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*, pp. 12922–12931, 2022.
- [24] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, “3d gaussian splatting for real-time radiance field rendering,” *ACM Transactions on Graphics*, July 2023.
- [25] 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.
- [26] 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.
- [27] 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.
- [28] Z. Jia, B. Wang, and C. Chen, “Drone-nerf: Efficient nerf based 3d scene reconstruction for large-scale drone survey,” *Image and Vision Computing*, vol. 143, p. 104920, 2024.
- [29] T. Suzuki, “Federated learning for large-scale scene modeling with neural radiance fields,” *arXiv preprint arXiv:2309.06030*, 2023.
- [30] S. Agarwal, Y. Furukawa, N. Snavely, I. Simon, B. Curless, S. M. Seitz, and R. Szeliski, “Building rome in a day,” *Communications of the ACM*, vol. 54, no. 10, pp. 105–112, 2011.
- [31] C. Früh and A. Zakhor, “An automated method for large-scale, ground-based city model acquisition,” *International Journal of Computer Vision*, vol. 60, pp. 5–24, 2004.
- [32] J. L. Schonberger and J.-M. Frahm, “Structure-from-motion revisited,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4104–4113, 2016.
- [33] 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*, pp. 8248–8258, 2022.

- [34] 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*, pp. 8296–8306, 2023.
- [35] Y. Deng, F. Lyu, J. Ren, Y. Zhang, Y. Zhou, Y. Zhang, and Y. Yang, "Share: Shaping data distribution at edge for communication-efficient hierarchical federated learning," in *2021 IEEE 41st International Conference on Distributed Computing Systems (ICDCS)*, pp. 24–34, IEEE, 2021.
- [36] Z. Li, Y. He, H. Yu, J. Kang, X. Li, Z. Xu, and D. Niyato, "Data heterogeneity-robust federated learning via group client selection in industrial iot," *IEEE Internet of Things Journal*, vol. 9, no. 18, pp. 17844–17857, 2022.
- [37] N. Mhaisen, A. A. Abdellatif, A. Mohamed, A. Erbad, and M. Guizani, "Optimal user-edge assignment in hierarchical federated learning based on statistical properties and network topology constraints," *IEEE Transactions on Network Science and Engineering*, vol. 9, no. 1, pp. 55–66, 2021.
- [38] Z. Zhong, W. Bao, J. Wang, X. Zhu, and X. Zhang, "Flee: A hierarchical federated learning framework for distributed deep neural network over cloud, edge, and end device," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 13, no. 5, pp. 1–24, 2022.
- [39] Y. Deng, J. Ren, C. Tang, F. Lyu, Y. Liu, and Y. Zhang, "A hierarchical knowledge transfer framework for heterogeneous federated learning," in *IEEE INFOCOM 2023-IEEE Conference on Computer Communications*, pp. 1–10, IEEE, 2023.
- [40] Z. Yang, S. Fu, W. Bao, D. Yuan, and A. Y. Zomaya, "Hierarchical federated learning with momentum acceleration in multi-tier networks," *IEEE Transactions on Parallel and Distributed Systems*, 2023.
- [41] E. Chen, S. Wang, and C. G. Brinton, "Taming subnet-drift in d2d-enabled fog learning: A hierarchical gradient tracking approach," in *IEEE INFOCOM 2024-IEEE conference on computer communications*, 2024. Accepted.