# *Scaffold*-GS: Structured 3D Gaussians for View-Adaptive Rendering

Tao Lu [1,3*]   Mulin Yu[1*]   Linning Xu[2]   Yuanbo Xiangli[4]
Limin Wang[1,3]   Dahua Lin[1,2]   Bo Dai[1]
[1]Shanghai Artificial Intelligence Laboratory, [2]The Chinese University of Hong Kong,
[3]Nanjing University, [4]Cornell University

Figure 1. *Scaffold*-**GS** represents the scene using a set of 3D Gaussians structured in a dual-layered hierarchy. Anchored on a sparse grid of initial points, a modest set of neural Gaussians are spawned from each anchor to *dynamically adapt to* various viewing angles and distances. Our method achieves rendering quality and speed comparable to 3D-GS with a more compact model (last row metrics: PSNR/storage size/FPS). Across multiple datasets, *Scaffold*-GS demonstrates more robustness in large outdoor scenes and intricate indoor environments with challenging observing views *e.g.* transparency, specularity, reflection, texture-less regions and fine-scale details.

## Abstract

*Neural rendering methods have significantly advanced photo-realistic 3D scene rendering in various academic and industrial applications. The recent 3D Gaussian Splatting method has achieved the state-of-the-art rendering quality and speed combining the benefits of both primitive-based representations and volumetric representations. However, it often leads to heavily redundant Gaussians that try to fit every training view, neglecting the underlying scene geometry. Consequently, the resulting model becomes less robust to significant view changes, texture-less area and lighting effects. We introduce Scaffold-GS, which uses anchor points to distribute local 3D Gaussians, and predicts their attributes on-the-fly based on viewing direction and distance within the view frustum. Anchor growing and pruning strategies are developed based on the importance of neural Gaussians to reliably improve the scene coverage. We show that our method effectively reduces redundant Gaussians while delivering high-quality rendering. We also demonstrates an enhanced capability to accommodate scenes with varying levels-of-detail and view-dependent ob-*

*servations, without sacrificing the rendering speed. Project page: https://city-super.github.io/scaffold-gs/.*

## 1. Introduction

Photo-realistic and real-time rendering of 3D scenes has always been a pivotal interest in both academic research and industrial domains, with applications spanning virtual reality [51], media generation [36], and large-scale scene visualization [43, 45, 49]. Traditional primitive-based representations like meshes and points [6, 26, 32, 55] are faster due to the use of rasterization techniques tailored for modern GPUs. However, they often yield low-quality renderings, exhibiting discontinuity and blurry artifacts. In contrast, volumetric representations and neural radiance fields utilize learning-based parametric models [3, 5, 30], hence can produce continuous rendering results with more details preserved. Nevertheless, they come with the cost of time-consuming stochastic sampling, leading to slower performance and potential noise.

In recent times, 3D Gaussian Splatting (3D-GS) [22] has

---

* denotes equal contribution.

1

achieved state-of-the-art rendering quality and speed. Initialized from point clouds derived from Structure from Motion (SfM) [42], this method optimizes a set of 3D Gaussians to represent the scene. It preserves the inherent continuity found in volumetric representations, whilst facilitating rapid rasterization by splatting 3D Gaussians onto 2D image planes.

While this approach offers several advantages, it tends to excessively expand Gaussian balls to accommodate every training view, thereby neglecting scene structure. This results in significant redundancy and limits its scalability, particularly in the context of complex large-scale scenes. Furthermore, view-dependent effects are baked into individual Gaussian parameters with little interpolation capabilities, making it less robust to substantial view changes and lighting effects.

We present *Scaffold*-GS, a Gaussian-based approach that utilizes anchor points to establish a hierarchical and region-aware 3D scene representation. We construct a sparse grid of anchor points initiated from SfM points. Each of these anchors tethers a set of neural Gaussians with learnable offsets, whose attributes (*i.e.* opacity, color, rotation, scale) are dynamically predicted based on the anchor feature and the viewing position. Unlike the vanilla 3D-GS which allows 3D Gaussians to freely drift and split, our strategy exploits scene structure to guide and constrain the distribution of 3D Gaussians, whilst allowing them to *locally* adaptive to varying viewing angles and distances. We further develop the corresponding growing and pruning operations for anchors to enhance the scene coverage.

Through extensive experiments, we show that our method delivers rendering quality on par with or even surpassing the original 3D-GS. At inference time, we limit the prediction of neural Gaussians to anchors within the view frustum, and filter out trivial neural Gaussians based on their opacity with a filtering step (*i.e.* learnable selector). As a result, our approach can render at a similar speed (around 100 FPS at 1K resolution) as the original 3D-GS with little computational overhead. Moreover, our storage requirements are significantly reduced as we only need to store anchor points and MLP predictors for each scene.

In conclusion, our contributions are: 1) Leveraging scene structure, we initiate *anchor points* from a sparse voxel grid to guide the distribution of local 3D Gaussians, forming a hierarchical and region-aware scene representation; 2) Within the view frustum, we predict *neural* Gaussians from each anchor *on-the-fly* to accommodate diverse viewing directions and distances, resulting in more robust novel view synthesis; 3) We develop a *more reliable* anchor growing and pruning strategy utilizing the predicted neural Gaussians for better scene coverage.

## 2. Related work

**MLP-based Neural Fields and Rendering.** Early neural fields typically adopt a multi-layer perceptron (MLP) as the global approximator of 3D scene geometry and appearance. They directly use spatial coordinates (and viewing direction) as input to the MLP and predict point-wise attribute, *e.g.* signed distance to scene surface (SDF) [33, 34, 46, 54], or density and color of that point [2, 30, 49]. Because of its volumetric nature and inductive bias of MLPs, this stream of methods achieves the SOTA performance in novel view synthesis. The major challenge of this scene representation is that the MLP need to be evaluated on a large number of sampled points along each camera ray. Consequently, rendering becomes extremely slow, with limited scalability towards complex and large-scale scenes. Despite several works have been proposed to accelerate or mitigate the intensive volumetric ray-marching, *e.g.* using proposal network [4], baking technique [11, 19], and surface rendering [41]. They either incorporated more MLPs or traded rendering quality for speed.

**Grid-based Neural Fields and Rendering.** This type of scene representations are usually based on a dense uniform grid of voxels. They have been greatly used in 3D shape and geometry modeling [12, 15, 21, 29, 35, 44, 57]. Some recent methods have also focused on faster training and inference of radiance field by exploiting spatial data structure to store scene features, which were interpolated and queried by sampled points during ray-marching. For instance, Plenoxel [13] adopted a sparse voxel grid to interpolate a continuous density field, and represented view-dependent visual effects with Spherical Harmonics. The idea of tensor factorization has been studied in multiple works [9, 10, 50, 52] to further reduce data redundancy and speed-up rendering. K-planes [14] used neural planes to parameterize a 3D scene, optionally with an additional temporal plane to accommodate dynamics. Several generative works [8, 40] also capitalized on triplane structure to model a 3D latent space for better geometry consistency. Instant-NGP [31] used a hash grid and achieved drastically faster feature query, enabling real-time rendering of neural radiance field. Although these approaches can produce high-quality results and are more efficient than global MLP representation, they still need to query many samples to render a pixel, and struggle to represent empty space effectively.

**Point-based Neural Fields and Rendering.** Point-based representations utilize the geometric primitive (*i.e.* point clouds) for scene rendering. A typical procedure is to rasterize an unstructured set of points using a fixed size, and exploits specialized modules on GPU and graphics APIs for rendering [7, 37, 38]. In spite of its fast speed and flexibil-

2

ity to solve topological changes, they usually suffer from holes and outliers that lead to artifacts in rendering. To alleviate the discontinuity issue, differentiable point-based rendering has been extensively studied to model objects geometry [16, 20, 27, 48, 55]. In particular, [48, 55] used differentiable surface splatting that treats point primitives as discs, ellipsoids or surfels that are larger than a pixel. [1, 24] augmented points with neural features and rendered using 2D CNNs. As a comparison, Point-NeRF [53] achieved high-quality novel view synthesis utilizing 3D volume rendering, along with region growing and point pruning during optimization. However, they resorted to volumetric ray-marching, hence hindered display rate. Notably, the recent work 3D-GS [22] employed anisotropic 3D Gaussians initialized from structure from motion (SfM) to represent 3D scenes, where a 3D Gaussian was optimized as a volume and projected to 2D to be rasterized as a primitive. Since it integrated pixel color using $\alpha$-blender, 3D-GS produced high-quality results with fine-scale detail, and rendered at real-time frame rate.

## 3. Methods

The original 3D-GS [22] optimizes Gaussians to reconstruct every training view, with heuristic splitting and pruning operations but in general neglects the underlying scene structure. This often leads to highly redundant Gaussians and makes the model less robust to novel viewing angles and distances. To address this issue, we propose a hierarchical 3D Gaussian scene representation that respects the scene geometric structure, with *anchor points* initialized from SfM to encode local scene information and spawn local *neural Gaussians*. The physical properties of neural Gaussians are decoded from the learned anchor features in a view-dependent manner *on-the-fly*. Fig. 2 illustrates our framework. We start with a brief background of 3D-GS then unfold our proposed method in details. Sec. 3.2.1 introduces how to initialize the scene with a regular sparse grid of anchor points from the irregular SfM point clouds. Sec. 3.2.2 explains how we predict neural Gaussians properties based on anchor points and view-dependent information. To make our method more robust to the noisy initialization, Sec. 3.3 introduces a neural Gaussian based "growing" and "pruning" operations to refine the anchor points. Sec. 3.4 elaborates training details.

### 3.1. Preliminaries

3D-GS [22] represents the scene with a set of anisotropic 3D Gaussians that inherit the differential properties of volumetric representation while be efficiently rendered via a tile-based rasterization.

Starting from a set of Structure-from-Motion (SfM) points, each point is designated as the position (mean) $\mu$

of a 3D Gaussian:

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

where $x$ is an arbitrary position within the 3D scene and $\Sigma$ denotes the covariance matrix of the 3D Gaussian. $\Sigma$ is formulated using a scaling matrix $S$ and rotation matrix $R$ to maintain its positive semi-definite:

$$\Sigma = RSS^T R^T, \qquad (2)$$

In addition to color $c$ modeled by Spherical harmonics, each 3D Gaussian is associated with an opacity $\alpha$ which is multiplied by $G(x)$ during the blending process.

Distinct from conventional volumetric representations, 3D-GS efficiently renders the scene via tile-based rasterization instead of resource-intensive ray-marching. The 3D Gaussian $G(x)$ are first transformed to 2D Gaussians $G'(x)$ on the image plane following the projection process as described in [58]. Then a tile-based rasterizer is designed to efficiently sort the 2D Gaussians and employ $\alpha$-blending:

$$C(x') = \sum_{i \in N} c_i \sigma_i \prod_{j=1}^{i-1}(1 - \sigma_j), \quad \sigma_i = \alpha_i G_i'(x'), \quad (3)$$

where $x'$ is the queried pixel position and N denotes the number of sorted 2D Gaussians associated with the queried pixel. Leveraging the differentiable rasterizer, all attributes of the 3D Gaussians are learnable and directly optimized end-to-end via training view reconstruction.

### 3.2. *Scaffold*-GS

#### 3.2.1 Anchor Point Initialization

Consistent with existing methods [22, 53], we use the sparse point cloud from COLMAP [39] as our initial input. We then voxelize the scene from this point cloud $\mathbf{P} \in \mathbb{R}^{M \times 3}$ as:

$$\mathbf{V} = \left\{ \left\lfloor \frac{\mathbf{P}}{\epsilon} \right\rceil \right\} \cdot \epsilon, \qquad (4)$$

where $\mathbf{V} \in \mathbb{R}^{N \times 3}$ denotes voxel centers, and $\epsilon$ is the voxel size. We then remove duplicate entries, denoted by $\{\cdot\}$ to reduce the redundancy and irregularity in $\mathbf{P}$.

The center of each voxel $v \in \mathbf{V}$ is treated as an anchor point, equipped with a local context feature $f_v \in \mathbb{R}^{32}$, a scaling factor $l_v \in \mathbb{R}^3$, and $k$ learnable offsets $\mathbf{O}_v \in \mathbb{R}^{k \times 3}$. In a slight abuse of terminology, we will denote the anchor point as $v$ in the following context. We further enhance $f_v$ to be multi-resolution and view-dependent. For each anchor $v$, we 1) create a features bank $\{f_v, f_{v_{\downarrow_1}}, f_{v_{\downarrow_2}}\}$, where $\downarrow_n$ denotes $f_v$ being down-sampled by $2^n$ factors; 2) blend the feature bank with view-dependent weights to form an integrated anchor feature $\hat{f}_v$. Specifically, Given a camera at
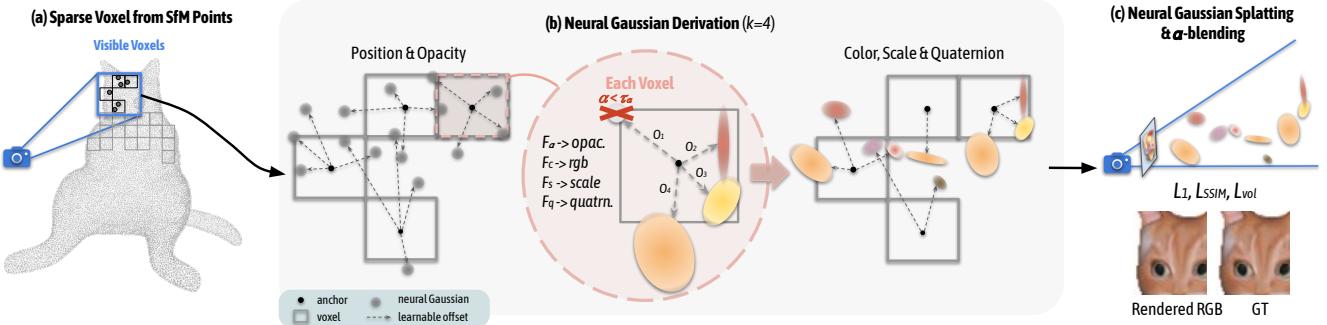
Figure 2. **Overview of *Scaffold*-GS.** (a) We start by forming a *sparse voxel grid* from SfM-derived points. An **anchor** associated with a learnable scale is placed at the center of each voxel, roughly sculpturing the scene occupancy. (b) Within a view frustum, $k$ **neural Gaussians** are spawned from each *visible anchor* with offsets $\{\mathcal{O}_k\}$. Their attributes, *i.e.* opacity, color, scale and quaternion are then decoded from the anchor feature, relative camera-anchor viewing direction and distance using $F_\alpha, F_c, F_s, F_q$ respectively. (c) Note that to alleviate redundancy and improve efficiency, only non-trivial neural Gussians (*i.e.* $\alpha \geq \tau_\alpha$) are rasterized following [22]. The rendered image is supervised via reconstruction ($\mathcal{L}_1$), structural similarity ($\mathcal{L}_{SSIM}$) and a volume regularization ($\mathcal{L}_{vol}$).

position $\mathbf{x}_c$ and an anchor at position $\mathbf{x}_v$, we calculate their relative distance and viewing direction with:

$$\delta_{vc} = \|\mathbf{x}_v - \mathbf{x}_c\|_2, \vec{\mathbf{d}}_{vc} = \frac{\mathbf{x}_v - \mathbf{x}_c}{\|\mathbf{x}_v - \mathbf{x}_c\|_2}, \quad (5)$$

then weighted sum the feature bank with weights predicted from a tiny MLP $F_w$:

$$\{w, w_1, w_2\} = \text{Softmax}(F_w(\delta_{vc}, \vec{\mathbf{d}}_{vc})), \quad (6)$$

$$\hat{f}_v = w \cdot f_v + w_1 \cdot f_{v_{\downarrow_1}} + w_2 \cdot f_{v_{\downarrow_2}}, \quad (7)$$

### 3.2.2 Neural Gaussian Derivation

In this section, we elaborate on how our approach derives neural Gaussians from anchor points. Unless specified otherwise, $F_*$ represents a particular MLP throughout the section. Moreover, we introduce two efficient pre-filtering strategies to reduce MLP overhead.

We parameterize a neural Gaussian with its position $\mu \in \mathbb{R}^3$, opacity $\alpha \in \mathbb{R}$, covariance-related quaternion $q \in \mathbb{R}^4$ and scaling $s \in \mathbb{R}^3$, and color $c \in \mathbb{R}^3$. As shown in Fig. 2(b), for each visible anchor point within the viewing frustum, we spawn $k$ neural Gaussians and predict their attributes. Specifically, given an anchor point located at $\mathbf{x}_v$, the positions of its neural Gaussians are calculated as:

$$\{\mu_0, ..., \mu_{k-1}\} = \mathbf{x}_v + \{\mathcal{O}_0, \ldots, \mathcal{O}_{k-1}\} \cdot l_v, \quad (8)$$

where $\{\mathcal{O}_0, \mathcal{O}_1, ..., \mathcal{O}_{k-1}\} \in \mathbb{R}^{k \times 3}$ are the learnable offsets and $l_v$ is the scaling factor associated with that anchor, as described in Sec. 3.2.1. The attributes of $k$ neural Gaussians are directly decoded from the anchor feature $\hat{f}_v$, the relative viewing distance $\delta_{vc}$ and direction $\vec{\mathbf{d}}_{vc}$ between the camera and the anchor point (Eq. 5) through individual MLPs,

denoted as $F_\alpha$, $F_c$, $F_q$ and $F_s$. Note that attributes are decoded in *one-pass*. For example, opacity values of neural Gaussians spawned from an anchor point are given by:

$$\{\alpha_0, ..., \alpha_{k-1}\} = F_\alpha(\hat{f}_v, \delta_{vc}, \vec{\mathbf{d}}_{vc}), \quad (9)$$

their colors $\{c_i\}$, quaternions $\{q_i\}$ and scales $\{s_i\}$ are similarly derived. Implementation details are in supplementary.

Note that the prediction of neural Gaussian attributes are *on-the-fly*, meaning that only anchors visible within the frustum are activated to spawn neural Gaussians. To make the rasterization more efficient, we only keep neural Gaussians whose opacity values are larger than a predefined threshold $\tau_\alpha$. This substantially cuts down the computational load and helps our method maintain a high rendering speed on-par with the original 3D-GS.

### 3.3. Anchor Points Refinement

**Growing Operation.** Since neural Gaussians are closely tied to their anchor points which are initialized from SfM points, their modeling power is limited to a local region, as has been pointed out in [22, 53]. This poses challenges to the initial placement of anchor points, especially in texture-less and less observed areas. We therefore propose an error-based anchor growing policy that grows new anchors where neural Gaussians find *significant*. To determine a *significant* area, we first spatially quantize the neural Gaussians by constructing voxels of size $\epsilon_g$. For each voxel, we compute the averaged gradients of the included neural Gaussians over $N$ training iterations, denoted as $\nabla_g$. Then, voxels with $\nabla_g > \tau_g$ is deemed as *significant*, where $\tau_g$ is a pre-defined threshold; and a new anchor point is thereby deployed at the center of that voxel if there was no anchor point established. Fig. 3 illustrates this growing operation. In practice, we quantize the space into multi-resolution voxel grid to al-
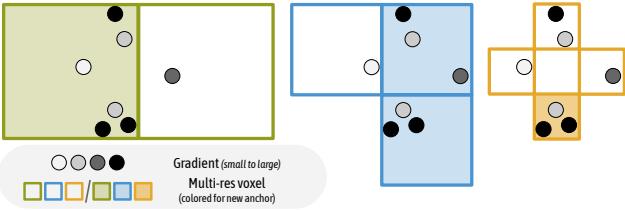
4

Figure 3. **Growing operation.** We develop an anchor growing policy guided by the gradients of the neural Gaussians. From left to right, we spatially quantize neural Gaussians into multi-resolution voxels ($m \in \{1, 2, 3\}$) of size $\{\epsilon_g^{(m)}\}$. New anchors are added to voxels with aggregated gradients larger than $\{\tau_g^{(m)}\}$.

low new anchors to be added at different granularity, where

$$\epsilon_g^{(m)} = \epsilon_g / 4^{m-1}, \quad \tau_g^{(m)} = \tau_g * 2^{m-1}, \quad (10)$$

where $m$ denotes the level of quantization. To further regulate the addition of new anchors, we apply a random elimination to these candidates. This cautious approach to adding points effectively curbs the rapid expansion of anchors.

**Pruning Operation**  To eliminate *trivial* anchors, we accumulate the opacity values of their associated neural Gaussians over $N$ training iterations. If an anchor fails to produce neural Gaussians with a satisfactory level of opacity, we then remove it from the scene.

### 3.4. Losses Design

We optimize the learnable parameters and MLPs with respect to the $\mathcal{L}_1$ loss over rendered pixel colors, with SSIM term [47] $\mathcal{L}_{SSIM}$ and volume regularization [28] $\mathcal{L}_{vol}$. The total supervision is given by:

$$\mathcal{L} = \mathcal{L}_1 + \lambda_{SSIM}\mathcal{L}_{SSIM} + \lambda_{vol}\mathcal{L}_{vol}, \quad (11)$$

where the volume regularization $\mathcal{L}_{vol}$ is:

$$\mathcal{L}_{vol} = \sum_{i=1}^{N_{ng}} \text{Prod}(s_i). \quad (12)$$

Here, $N_{ng}$ denotes the number of neural Gaussians in the scene and $\text{Prod}(\cdot)$ is the product of the values of a vector, *e.g.*, in our case the scale $s_i$ of each neural Gaussian. The volume regularization term encourages the neural Gaussians to be small with minimal overlapping.

## 4. Experiments

### 4.1. Experimental Setup

**Dataset and Metrics.**  We conducted a comprehensive evaluation across 27 scenes from publicly available datasets. Specifically, we tested our approach on

all available scenes tested in the 3D-GS [22], including seven scenes from Mip-NeRF360 [4], two scenes from Tanks&Temples [23], two scenes from DeepBlending [18] and synthetic Blender dataset [30]. We additionally evaluated on datasets with contents captured at multiple LODs to demonstrate our advantages in view-adaptive rendering. Six scenes from BungeeNeRF [49] and two scenes from VR-NeRF [51] are selected. The former provides multi-scale outdoor observations and the latter captures intricate indoor environments. Apart from the commonly used metrics (PSNR, SSIM [47], and LPIPS [56]), we additionally report the storage size (MB) and the rendering speed (FPS) for model compactness and performance efficiency. We provide the averaged metrics over all scenes of each dataset in the main paper and leave the full quantitative results on each scene in the supplementary.

**Baseline and Implementation.**  3D-GS [22] is selected as our main baseline for its established SOTA performance in novel view synthesis. Both 3D-GS and our method were trained for 30k iterations. We also record the results of Mip-NeRF360 [4], iNGP [31] and Plenoxels [13] as in [22].

For our method, we set $k = 10$ for all experiments. All the MLPs employed in our approach are 2-layer MLPs with ReLU activation; the dimensions of the hidden units are all 32. For anchor points refinement, we average gradients over $N = 100$ iterations, and by default use $\tau_g = 64\epsilon$. On intricate scenes and the ones with dominant texture-less regions, we use $\tau_g = 16\epsilon$. An anchor is pruned if the accumulated opacity of its neural Gaussians is less than $0.5$ at each round of refinement. The two loss weights $\lambda_{SSIM}$ and $\lambda_{vol}$ are set to $0.2$ and $0.001$ in our experiments. Please check the supplementary material for more details.

### 4.2. Results Analysis

Our evaluation was conducted on diverse datasets, ranging from synthetic object-level scenes, indoor and outdoor environments, to large-scale urban scenes and landscapes. A variety of improvements can be observed especially on challenging cases, such as texture-less area, insufficient observations, fine-scale details and view-dependent light effects. See Fig. 1 and Fig. 4 for examples.

**Comparisons.**  In assessing the quality of our approach, we compared with 3D-GS [22], Mip-NeRF360 [4], iNGP [31] and Plenoxels [13] on real-world datasets. Qualitative results are presented in Tab. 1. The quality metrics for Mip-NeRF360, iNGP and Plenoxels align with those reported in the 3D-GS study. It can be noticed that our approach achieves comparable results with the SOTA algorithms on Mip-NeRF360 dataset, and surpassed the SOTA on Tanks&Temples and DeepBlending, which captures more challenging environments with the presence

Table 1. **Quantitative comparison to previous methods on real-world datasets.** Competing metrics are extracted from respective papers.

| Dataset | Mip-NeRF360 | | | Tanks&Temples | | | Deep Blending | | |
|---|---|---|---|---|---|---|---|---|---|
| Method \ Metrics | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
| **3D-GS** [22] | 28.69 | 0.870 | 0.182 | 23.14 | 0.841 | 0.183 | 29.41 | 0.903 | 0.243 |
| **Mip-NeRF360** [4] | 29.23 | 0.844 | 0.207 | 22.22 | 0.759 | 0.257 | 29.40 | 0.901 | 0.245 |
| **iNPG** [31] | 26.43 | 0.725 | 0.339 | 21.72 | 0.723 | 0.330 | 23.62 | 0.797 | 0.423 |
| **Plenoxels** [13] | 23.62 | 0.670 | 0.443 | 21.08 | 0.719 | 0.379 | 23.06 | 0.795 | 0.510 |
| **Ours** | 28.84 | 0.848 | 0.220 | 23.96 | 0.853 | 0.177 | 30.21 | 0.906 | 0.254 |



Figure 4. **Qualitative comparison of *Scaffold*-GS and 3D-GS [22] across diverse datasets [4, 17, 23, 51].** Patches that highlight the visual differences are emphasized with arrows and green & yellow insets for clearer visibility. Our approach constantly outperforms 3D-GS on these scenes, with evident advantages in challenging scenarios, *e.g.* thin geometry and fine-scale details (MIP360-ROOM(a), MIP360-COUNTER(a)), texture-less regions (DB-DRJOHNSON, DB-PLAYROOM), light effects (MIP360-COUNTER(b), DB-DRJOHNSON), insufficient observations (TANDT-TRAIN, VR-KITCHEN). It can also be observed (*e.g.* VR-APARTMENT) that our model is superior in representing contents at varying scales and viewing distances.

Table 2. **Performance comparison.** Rendering FPS and storage size are reported. Storage size reduction ratio is indicated by (↓). Rendering speed of both methods are measured on our machine.

| Dataset | Mip-NeRF360 | | Tanks&Temples | | Deep Blending | |
|---|---|---|---|---|---|---|
| | FPS | Mem (MB) | FPS | Mem (MB) | FPS | Mem (MB) |
| **3D-GS** | 97 | 693 | **123** | 411 | 109 | 676 |
| **Ours** | **102** | **156** (4.4× ↓) | 110 | **87** (4.7× ↓) | **139** | **66** (10.2× ↓) |

of *e.g.* changing lighting, texture-less regions and reflections. In terms of efficiency, we evaluated rendering speed and storage size of our method and 3D-GS, as shown in

Tab. 2. Our method achieved real-time rendering while using less storage, indicating that our model is more compact than 3D-GS without sacrificing rendering quality and speed. Additionally, akin to prior grid-based methods, our approach converged faster than 3D-GS. See supplementary material for more analysis.

We also examined our method on the synthetic Blender dataset, which provides an exhaustive set of views capturing objects at $360°$. A good set of initial SfM points is not readily available in this dataset, thus we start from $100k$ grid points and learn to grow and prune points with our anchor refinement operations. After 30k iterations, we used the re-

Table 3. **Qualitative comparison.** Our method is able to handle large-scale scenes (*e.g.* BUNGEENERF) with light-weight representation. Our method shows consistent compactness and effectiveness in complex lighting conditions and synthetic scenes.

| Dataset | BungeeNeRF | | VR-NeRF | | Synthetic Blender | |
|---------|------|----------|------|----------|------|----------|
| | PSNR | Mem (MB) | PSNR | Mem (MB) | PSNR | Mem (MB) |
| **3D-GS** | 24.89 | 1606 | 28.94 | 263 | 33.32 | 53 |
| **Ours** | **27.01** | **203** (7.9× ↓) | **29.24** | **69** (3.8× ↓) | **33.68** | **14** (3.8× ↓) |



Figure 5. **Comparison on multi-scale scenes (w/ zoom-in cases).** We showcase the rendering outcomes at an unseen closer scale on the AMSTERDAM scene from BungeeNeRF. Our method smoothly extrapolates to new viewing distances using refined neural Gaussian properties, remedying the needle-like artifacts of original 3D-GS caused by fixed Gaussian scaling values.

mained points as initialized anchors and re-run our framework. The PSNR score and storage size compared with 3D-GS are presented in Tab. 3. Fig. 1 also demonstrates that our method can achieve better visual quality with more reliable geometry and texture details.

**Multi-scale Scene Contents.** We examined our model's capability in handling multi-scale scene details on the BungeeNeRF and VR-NeRF datasets. As shown in Tab. 3, our method achieved superior quality whilst using fewer storage size to store the model compared to 3D-GS [22]. As illustrated in Fig. 4 and Fig. 5, our method was superior in accommodating varying levels of detail in the scene. In contrast, images rendered from 3D-GS often suffered from noticeable blurry and needle-shaped artifacts. This is likely because that Gaussian attributes are optimized to overfit multi-scale training views, creating excessive Gaussians that work for each observing distance. However, it can easily lead to ambiguity and uncertainty when synthesizing novel views, since it lacks the ability to reason about viewing angle and distance. On contrary, our method efficiently encoded local structures into compact neural features, enhancing both rendering quality and convergence speed. Details are provided in the supplementary material.



Figure 6. **Anchor feature clustering.** We cluster anchor features (DB-PLAYROOM) into 3 clusters using K-means [25] and visualize the result. The clustered features show clues of scene contents, *e.g.* the banister, stroller, desk and monitor can be clearly identified. Anchors on the wall and floor are also respectively grouped together. This shows that our approach improves the interpretability of 3D-GS model, and has the potential to be scaled-up on much larger scenes exploiting reusable features.

**Feature Analysis.** We further perform an analysis of the learnable anchor features and the selector mechanism. As depicted in Fig. 6, the clustered pattern suggests that the compact anchor feature spaces adeptly capture regions with similar visual attributes and geometries, as evidenced by their proximity in the encoded feature space.

**View Adaptability.** To support that our neural Gaussians are view-adaptive, we explore how the values of attributes change when the same Gaussian is observed from different positions. Fig. 7 demonstrates a varying distribution of attributes intensity at different viewing positions, while maintaining a degree of local continuity. This accounts for the superior *view adaptability* of our method compared to the static attributes of 3D-GS, as well as its enhanced generalizability to novel views.

**Selection Process by Opacity.** We examine the decoded opacity from neural Gaussians and our opacity-based selection process (Fig. 2(b)) from two aspects. First, without the anchor point refinement module, we filter neural Gaussian using their decoded opacity values to extract geometry from a random point cloud. Fig. 8 demonstrates that the remained neural Gaussians effectively reconstruct the coarse structure of the bulldozer model from random points, highlighting its capability for implicit geometry modeling under mainly rendering-based supervision. We found this is conceptually similar to the proposal network utilized in [4], serving as the geometry proxy estimator for efficient sampling.
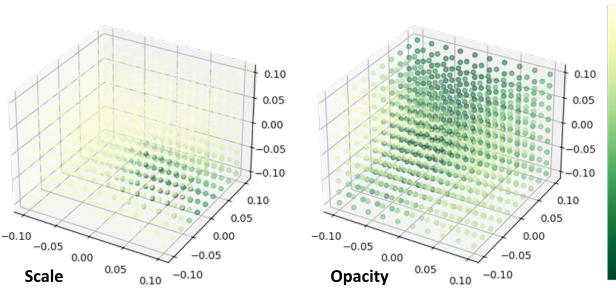
Second, we apply different $k$ values in our method.

Figure 7. **View-adaptive neural Gaussian attributes.** We visualize the decoded attributes of a *single* neural Gaussian observed at different positions. Each point corresponds to a viewpoint in space. The color of the point denotes the intensity of attributes decoded for this view (left: $F_s \rightarrow s_i$; right: $F_\alpha \rightarrow \alpha_i$). This pattern indicates that attributes of a neural Gaussian adapt to viewpoint changing, while exhibiting a certain degree of local continuity.
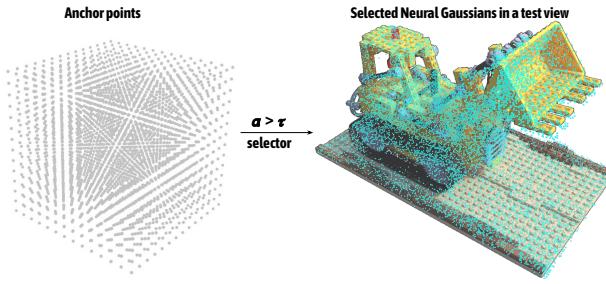


Figure 8. **Geometry culling via selector.** (Left) Anchor points from randomly initialized points; (Right) Activated neural Gaussians derived from each anchor under the current view. In synthetic Blender scenes, with all 3D Gaussians visible in the viewing frustum, our opacity filtering functions similar to a geometry proxy estimator, excluding unoccupied regions before rasterization.
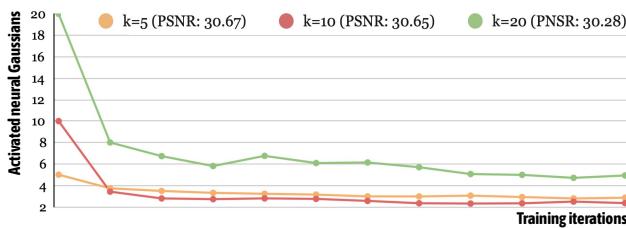


Figure 9. **Learning with different $k$ values.** Despite varying initial $k$ values under different hyper-parameter settings, they converge to activate a similar number of neural Gaussians with comparable rendering fidelity.

Fig. 9 shows that regardless of the $k$ value, the final number of activated neural Gaussians converges to a similar amount through the training, indicating *Scaffold*-GS's preference to select a collection of non-redundant Gaussians that are sufficient to represent the scene.

Table 4. **Effects of filtering**. FILTER 1 refers to selecting anchors by view frustum and FILTER 2 refers to the opacity-based selection process. The filtering method has no notable impact on fidelity, but greatly affects inference speed.

| Scene | DB-PLAYROOM | | DB-DRJOHNSON | |
|---|---|---|---|---|
| | PSNR | FPS | PSNR | FPS |
| NO FILTERS | 30.4 | 84 | 29.7 | 79 |
| FILTER 1 | 30.3 | 118 | 29.6 | 100 |
| FILTER 2 | 30.6 | 109 | 29.7 | 104 |
| FULL | 30.62 | 150 | 29.8 | 129 |

Table 5. **Anchor refinement**. The growing operation is essential for fidelity since it improves the poor initialization. The pruning operation controls the increasing of storage size and optimizes the quality of remained anchors.

| Scene | DB-PLAYROOM | | DB-DRJOHNSON | |
|---|---|---|---|---|
| | PSNR | Mem (MB) | PSNR | Mem (MB) |
| NONE | 28.45 | 24 | 28.81 | 12 |
| W/ PRUNING | 29.12 | 23 | 28.51 | 12 |
| W/ GROWING | 30.54 | 71 | 29.75 | 76 |
| FULL | 30.62 | 63 | 29.80 | 68 |

### 4.3. Ablation Studies

**Efficacy of Filtering Strategies.** We evaluated our filtering strategies (Sec. 3.2.2), which we found crucial for speeding up our method. As Tab. 4 shows, while these strategies had no notable effect on fidelity, they significantly enhanced inference speed. However, there was a risk of masking pertinent neural Gaussians, which we aim to address in future works.

**Efficacy of Anchor Points Refinement Policy.** We evaluated our growing and pruning operations described in Sec. 3.3. Tab. 5 shows the results of disabling each operation in isolation and maintaining the rest of the method. We found that the addition operation is crucial for accurately reconstructing details and texture-less areas, while the pruning operation plays an important role in eliminating trivial Gaussians and maintaining the efficiency of our approach.

### 4.4. Discussions and Limitations

Through our experiments, we found that the initial points play a crucial role for high-fidelity results. Initializing our framework from SfM point clouds is a swift and viable solution, considering these point clouds usually arise as a byproduct of image calibration processes. However, this approach may be suboptimal for scenarios dominated by large texture-less regions. Despite our anchor point refinement strategy can remedy this issue to some extent, it still suffers from extremely sparse points. We expect that our algorithm will progressively improve as the field advances, yielding

more accurate results. Further details are discussed in the supplementary material.

## 5. Conclusion

In this work, we introduce *Scaffold*-GS, a novel 3D neural scene representation for efficient view-adaptive rendering. The core of *Scaffold*-GS lies in its structural arrangement of 3D Gaussians guided by anchor points from SfM, whose attributes are on-the-fly decoded from view-dependent MLPs. We show that our approach leverages a much more compact set of Gaussians to achieve comparable or even better results than the SOTA algorithms. The advantage of our *view-adaptive* neural Gaussians is particularly evident in challenging cases where 3D-GS usually fails. We further show that our anchor points encode local features in a meaningful way that exhibits semantic patterns to some degree, suggesting its potential applicability in a range of versatile tasks such as large-scale modeling, manipulation and interpretation in the future.

## References

[1] Kara-Ali Aliev, Dmitry Ulyanov, and Victor S. Lempitsky. Neural point-based graphics. In *European Conference on Computer Vision*, 2019. 3

[2] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5835–5844, 2021. 2

[3] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5855–5864, 2021. 1

[4] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. *CVPR*, 2022. 2, 5, 6, 7, 3

[5] 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 (ICCV)*, pages 19697–19705, 2023. 1

[6] Mario Botsch, Alexander Hornung, Matthias Zwicker, and Leif Kobbelt. High-quality surface splatting on today's gpus. In *Proceedings Eurographics/IEEE VGTC Symposium Point-Based Graphics, 2005.*, pages 17–141. IEEE, 2005. 1

[7] Mario Botsch, Alexander Sorkine-Hornung, Matthias Zwicker, and Leif P. Kobbelt. High-quality surface splatting on today's gpus. *Proceedings Eurographics/IEEE VGTC Symposium Point-Based Graphics, 2005.*, pages 17–141, 2005. 2

[8] Eric Chan, Connor Z. Lin, Matthew Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J. Guibas, Jonathan Tremblay, S. Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3d generative adversarial networks. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16102–16112, 2021. 2

[9] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. *ArXiv*, abs/2203.09517, 2022. 2

[10] Anpei Chen, Zexiang Xu, Xinyue Wei, Siyu Tang, Hao Su, and Andreas Geiger. Factor fields: A unified framework for neural fields and beyond. *ArXiv*, abs/2302.01226, 2023. 2

[11] Zhiqin Chen, Thomas Funkhouser, Peter Hedman, and Andrea Tagliasacchi. Mobilenerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. In *The Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2

[12] Christopher Bongsoo Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. *ArXiv*, abs/1604.00449, 2016. 2

[13] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022. 2, 5, 6, 3

[14] Sara Fridovich-Keil, Giacomo Meanti, Frederik Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12479–12488, 2023. 2

[15] Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas A. Funkhouser. Local deep implicit functions for 3d shape. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4856–4865, 2019. 2

[16] Markus Gross and Hanspeter Pfister. *Point-based graphics*. Elsevier, 2011. 3

[17] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. Deep blending for free-viewpoint image-based rendering. 37(6):257:1–257:15, 2018. 6, 3

[18] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. Deep blending for free-viewpoint image-based rendering. *ACM Transactions on Graphics (ToG)*, 37(6):1–15, 2018. 5, 2, 3

[19] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul E. Debevec. Baking neural radiance fields for real-time view synthesis. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5855–5864, 2021. 2

[20] Eldar Insafutdinov and Alexey Dosovitskiy. Unsupervised learning of shape and pose with differentiable point clouds. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 3

[21] Abhishek Kar, Christian Häne, and Jitendra Malik. Learning a multi-view stereo machine. *ArXiv*, abs/1708.05375, 2017. 2

[22] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time

radiance field rendering. *ACM Transactions on Graphics*, 42 (4), 2023. 1, 3, 4, 5, 6, 7, 2

[23] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics*, 36(4), 2017. 5, 6, 2, 3

[24] Georgios Kopanas, Julien Philip, Thomas Leimkühler, and George Drettakis. Point-based neural rendering with per-view optimization. *Computer Graphics Forum*, 40, 2021. 3

[25] K Krishna and M Narasimha Murty. Genetic k-means algorithm. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 29(3):433–439, 1999. 7

[26] Christoph Lassner and Michael Zollhofer. Pulsar: Efficient sphere-based neural rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1440–1449, 2021. 1

[27] Chen-Hsuan Lin, Chen Kong, and Simon Lucey. Learning efficient point cloud generation for dense 3d object reconstruction. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2018. 3

[28] Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhoefer, Yaser Sheikh, and Jason Saragih. Mixture of volumetric primitives for efficient neural rendering. *ACM Transactions on Graphics (ToG)*, 40(4):1–13, 2021. 5

[29] Lars M. Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4455–4465, 2018. 2

[30] 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, 5, 3

[31] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022. 2, 5, 6, 3

[32] Jacob Munkberg, Jon Hasselgren, Tianchang Shen, Jun Gao, Wenzheng Chen, Alex Evans, Thomas Müller, and Sanja Fidler. Extracting Triangular 3D Models, Materials, and Lighting From Images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8280–8290, 2022. 1

[33] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5569–5579, 2021. 2

[34] Jeong Joon Park, Peter R. Florence, Julian Straub, Richard A. Newcombe, and S. Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 165–174, 2019. 2

[35] Songyou Peng, Michael Niemeyer, Lars M. Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. *ArXiv*, abs/2003.04618, 2020. 2

[36] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*, 2023. 1

[37] Liu Ren, Hanspeter Pfister, and Matthias Zwicker. Object space ewa surface splatting: A hardware accelerated approach to high quality point rendering. *Computer Graphics Forum*, 21, 2002. 2

[38] Miguel Sainz and Renato Pajarola. Point-based rendering techniques. *Computers & Graphics*, 28(6):869–879, 2004. 2

[39] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 3, 2

[40] Jessica Shue, Eric Chan, Ryan Po, Zachary Ankner, Jiajun Wu, and Gordon Wetzstein. 3d neural field generation using triplane diffusion. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20875–20886, 2022. 2

[41] Vincent Sitzmann, Semon Rezchikov, William T. Freeman, Joshua B. Tenenbaum, and Frédo Durand. Light field networks: Neural scene representations with single-evaluation rendering. In *Neural Information Processing Systems*, 2021. 2

[42] Noah Snavely, Steven M. Seitz, and Richard Szeliski. *Photo Tourism: Exploring Photo Collections in 3D*. Association for Computing Machinery, New York, NY, USA, 1 edition, 2023. 2

[43] Matthew Tancik, Vincent Casser, Xinchen Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. 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

[44] Shubham Tulsiani, Tinghui Zhou, Alyosha A. Efros, and Jitendra Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 209–217, 2017. 2

[45] 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 (CVPR)*, pages 12922–12931, 2022. 1

[46] 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

[47] Zhou 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*, 13(4): 600–612, 2004. 5

[48] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. Synsin: End-to-end view synthesis from a single image. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7465–7475, 2019. 3

[49] Yuanbo Xiangli, Linning 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 *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part XXXII*, pages 106–122. Springer, 2022. 1, 2, 5, 3

[50] Yuanbo Xiangli, Linning Xu, Xingang Pan, Nanxuan Zhao, Bo Dai, and Dahua Lin. Assetfield: Assets mining and reconfiguration in ground feature plane representation. *ArXiv*, abs/2303.13953, 2023. 2

[51] Linning Xu, Vasu Agrawal, William Laney, Tony Garcia, Aayush Bansal, Changil Kim, Samuel Rota Bulò, Lorenzo Porzi, Peter Kontschieder, Aljaž Božič, Dahua Lin, Michael Zollhöfer, and Christian Richardt. VR-NeRF: High-fidelity virtualized walkable spaces. In *SIGGRAPH Asia Conference Proceedings*, 2023. 1, 5, 6, 2, 3

[52] Linning 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

[53] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-nerf: Point-based neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5438–5448, 2022. 3, 4

[54] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021. 2

[55] Wang Yifan, Felice Serena, Shihao Wu, Cengiz Öztireli, and Olga Sorkine-Hornung. Differentiable surface splatting for point-based geometry processing. *ACM Transactions on Graphics (TOG)*, 38(6):1–14, 2019. 1, 3

[56] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 5

[57] Xi Zhao, Ruizhen Hu, Haisong Liu, Taku Komura, and Xinyu Yang. Localization and completion for 3d object interactions. *IEEE Transactions on Visualization and Computer Graphics*, 26(8):2634–2644, 2019. 2

[58] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Ewa volume splatting. In *Proceedings Visualization, 2001. VIS'01.*, pages 29–538. IEEE, 2001. 3

# *Scaffold*-GS: Structured 3D Gaussians for View-Adaptive Rendering

## Supplementary Material

## 6. Overview

This supplementary is organized as follows: (1) In the first section, we elaborate implementation details of our *Scaffold*-GS, including anchor point feature enhancement (Sec.3.2.1), structure of MLPs (Sec.3.2.2) and anchor point refinement strategies (Sec.3.3); (2) The second part describes our dataset preparation steps. We then show additional experimental results and analysis based on our training observations.

## 7. Implementation details.

**Feature Bank.** To enhance the view-adaptability, we update the anchor feature through a view-dependent encoding. Following calculating the relative distance $\delta_{vc}$ and viewing direction $\vec{\mathbf{d}}_{vc}$ of a camera and an anchor, we predict a weight vector $w \in \mathbb{R}^3$ as follows:

$$(w, w_1, w_2) = \text{Softmax}(F_w(\delta_{vc}, \vec{\mathbf{d}}_{vc})), \qquad (13)$$

where $F_w$ is a tiny MLP that serves as a view encoding function. We then encode the view direction information to the anchor feature $f_v$ by compositing a feature bank containing information with different resolutions as follows:

$$\hat{f}_v = w \cdot f_v + w_1 \cdot f_{v_{\downarrow_1}} + w_2 \cdot f_{v_{\downarrow_2}}, \qquad (14)$$

In practice, we implement the feature bank via slicing and repeating, as illustrated in Fig. 10. We found this slicing and mixture operation improves *Scaffold*-GS's ability to capture different scene granularity. The distribution of feature bank's weights is illustrated in Fig. 11.
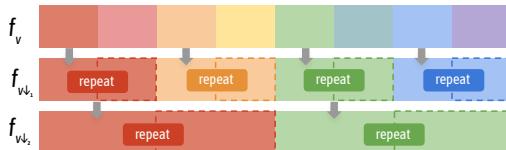


Figure 10. **Generation of Feature Bank.** We expand the anchor feature $f$ into a set of *multi-resolution* features $\{f_v, f_{v_{\downarrow_1}}, f_{v_{\downarrow_2}}\}$ via slicing and repeating. This operation improves *Scaffold*-GS's ability to capture different scene granularity.

**MLPs as feature decoders.** The core MLPs include the opacity MLP $F_\alpha$, the color MLP $F_c$ and the covariance MLP $F_s$ and $F_q$. All of these $F_*$ are implemented in a LINEAR → RELU → LINEAR style with the hidden dimension of 32, as illustrated in Fig. 12. Each branch's output is activated with a head layer.
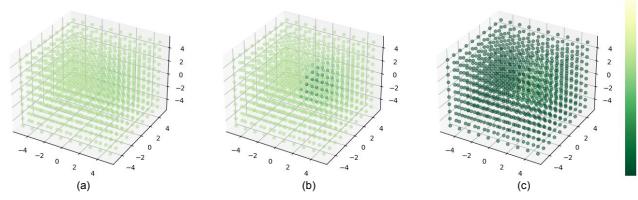


Figure 11. **View-based feature bank's weight distribution.** (a), (b) and (c) denote the predicted weights $\{w_2, w_1, w\}$ for $\{f_{v_{\downarrow_2}}, f_{v_{\downarrow_1}}, f_v\}$ from a group of uniformly distributed viewpoints. Light color denotes larger weights. For this anchor, finer features are more activated at center view positions. The patterns exhibit the ability to capture different scene granularities based on view direction and distance.
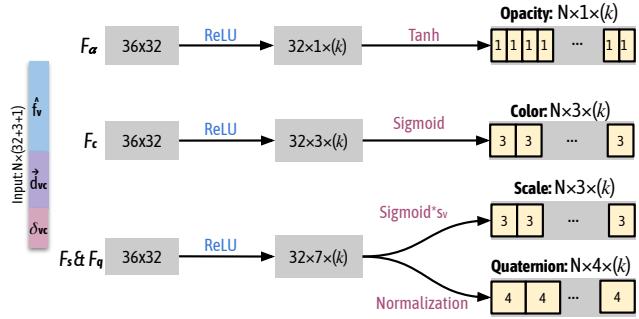


Figure 12. **MLP Structures.** For each anchor point, we use small MLPs ($F_\alpha, F_c, F_s, F_q$) to predict attributes (opacity, color, scale and quaternion) of $k$ neural Gaussians. The input to MLPs are anchor feature $\hat{f}_v$, relative viewing direction $\vec{\mathbf{d}}_{vc}$ and distance $\delta_{vc}$ between the camera and anchor point.

- For *opacity*, the output is activated by $\text{Tanh}$, where value 0 serves as a natural threshold for selecting valid samples and the final valid values can cover the full range of [0,1].
- For *color*, we activate the output with $\text{Sigmoid}$ function:

$$\{c_0, ..., c_{k-1}\} = \text{Sigmoid}(F_c), \qquad (15)$$

which constrains the color into a range of (0,1).
- For *rotation*, we follow 3D-GS [22] and activate it with a normalization to obtain a valid quaternion.
- For *scaling*, we adjust the base scaling $s_v$ of each anchor with the MLP output as follows:

$$\{s_0, ..., s_{k-1}\} = \text{Sigmoid}(F_s) \cdot s_v, \qquad (16)$$

**Voxel Size.** The voxel size $\epsilon$ sets the finest anchor resolution. We employ two strategies: 1) Use the median of the nearest-neighbor distances among all initial points: $\epsilon$
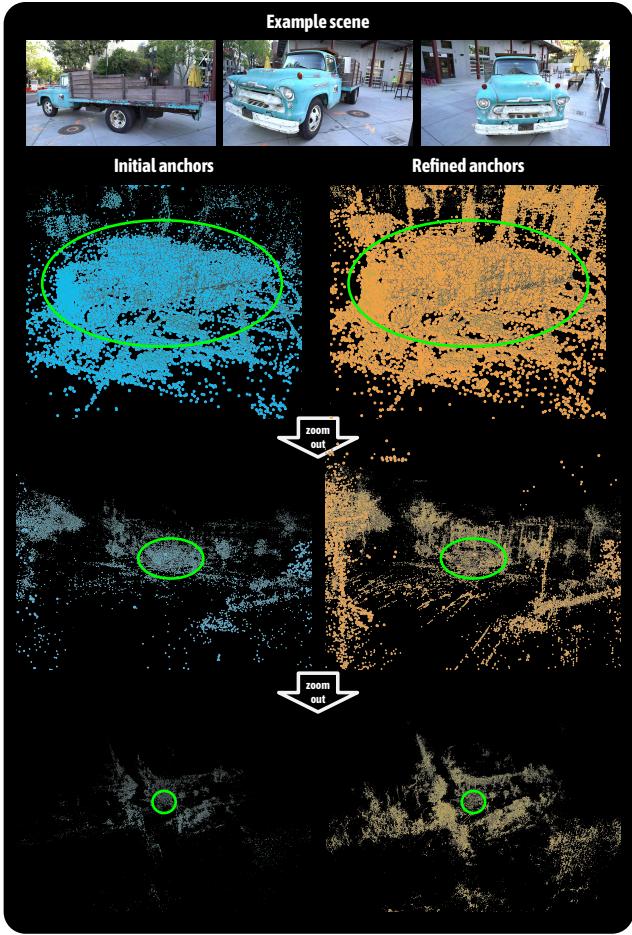
Figure 13. **Anchor Refinement.** We visualize the initial and refined anchor points on the truck scene [23]. The truck is highlighted by the circle. Note that the refined points effectively covers surrounding regions and fine-scale structures, leaning to more complete and detailed scene renderings.

is adapted to point cloud density, yielding denser anchors with enhanced rendering quality but might introduce more computational overhead; 2) Set $\epsilon$ manually to either 0.005 or 0.01: this is effective in most scenarios but might lead to missing details in texture-less regions. We found these two strategies adequately accommodate various scene complexities in our experiments.

**Anchor Refinement.** As briefly discussed in the main paper, the voxelization process suggests that our method may behave sensitive to initial SfM results. We illustrate the effect of the anchor refinement process in Fig. 13, where new anchors enhance scene details and fill gaps in large texture-less regions and less observed areas.

Table 6. **SSIM scores for Mip-NeRF360 [4] scenes.**

| Method | Scenes | bicycle | garden | stump | room | counter | kitchen | bonsai |
|---|---|---|---|---|---|---|---|---|
| **3D-GS [22]** | | **0.771** | **0.868** | 0.775 | 0.914 | 0.905 | 0.922 | 0.938 |
| **Mip-NeRF360 [4]** | | 0.685 | 0.813 | 0.744 | 0.913 | 0.894 | 0.920 | 0.941 |
| **iNPG [31]** | | 0.491 | 0.649 | 0.574 | 0.855 | 0.798 | 0.818 | 0.890 |
| **Plenoxels [13]** | | 0.496 | 0.6063 | 0.523 | 0.8417 | 0.759 | 0.648 | 0.814 |
| **Ours** | | 0.705 | 0.842 | **0.784** | **0.925** | **0.914** | **0.928** | **0.946** |

Table 7. **PSNR scores for Mip-NeRF360 [4] scenes.**

| Method | Scenes | bicycle | garden | stump | room | counter | kitchen | bonsai |
|---|---|---|---|---|---|---|---|---|
| **3D-GS [22]** | | **25.25** | **27.41** | **26.55** | 30.63 | 28.70 | 30.32 | 31.98 |
| **Mip-NeRF360 [4]** | | 24.37 | 26.98 | 26.40 | 31.63 | **29.55** | **32.23** | **33.46** |
| **iNPG [31]** | | 22.19 | 24.60 | 23.63 | 29.27 | 26.44 | 28.55 | 30.34 |
| **Plenoxels [13]** | | 21.91 | 23.49 | 20.66 | 27.59 | 23.62 | 23.42 | 24.67 |
| **Ours** | | 24.50 | 27.17 | 26.27 | **31.93** | 29.34 | 31.30 | 32.70 |

Table 8. **LPIPS scores for Mip-NeRF360 [4] scenes.**

| Method | Scenes | bicycle | garden | stump | room | counter | kitchen | bonsai |
|---|---|---|---|---|---|---|---|---|
| **3D-GS [22]** | | **0.205** | **0.103** | **0.210** | 0.220 | 0.204 | 0.129 | 0.205 |
| **Mip-NeRF360 [4]** | | 0.301 | 0.170 | 0.261 | 0.211 | 0.204 | 0.127 | **0.176** |
| **iNPG [31]** | | 0.487 | 0.312 | 0.450 | 0.301 | 0.342 | 0.254 | 0.227 |
| **Plenoxels [13]** | | 0.506 | 0.3864 | 0.503 | 0.4186 | 0.441 | 0.447 | 0.398 |
| **Ours** | | 0.306 | 0.146 | 0.284 | **0.202** | **0.191** | **0.126** | 0.185 |

Table 9. **Storage size (MB) for Mip-NeRF360 [4] scenes.**

| Method | Scenes | bicycle | garden | stump | room | counter | kitchen | bonsai |
|---|---|---|---|---|---|---|---|---|
| **3D-GS [22]** | | 1291 | 1268 | 1034 | 327 | 261 | 414 | 281 |
| **Ours** | | 248 | 271 | 493 | 133 | 194 | 173 | 258 |

Table 10. **SSIM scores for Tanks&Temples [23] and Deep Blending [18] scenes.**

| Method | Scenes | Truck | Train | Dr Johnson | Playroom |
|---|---|---|---|---|---|
| **3D-GS [22]** | | 0.879 | 0.802 | 0.899 | **0.906** |
| **Mip-NeRF360 [4]** | | 0.857 | 0.660 | 0.901 | 0.900 |
| **iNPG [31]** | | 0.779 | 0.666 | 0.839 | 0.754 |
| **Plenoxels [13]** | | 0.774 | 0.663 | 0.787 | 0.802 |
| **Ours** | | **0.883** | **0.822** | **0.907** | 0.904 |

## 8. Experiments and Results

**Additional Data Preprocessing.** We used COLMAP [39] to estimate camera poses and generate SfM points for VR-NeRF [51] and BungeeNeRF [49] datasets. Both two datasets are challenging in terms of varying levels of details presented in the captures. The VR-NeRF dataset was tested using its eye-level subset with 3 cameras. For all other datasets, we adhered to the original 3D-GS [22] method, sourcing them from public resources.

**Per-scene Results.** Here we list the error metrics used in our evaluation in Sec.4 across all considered methods and scenes, as shown in Tab. 6-17.

**Training Process Analysis.** Figure 14 illustrates the variations in PSNR during the training process for both training and testing views. Our method demonstrates quicker
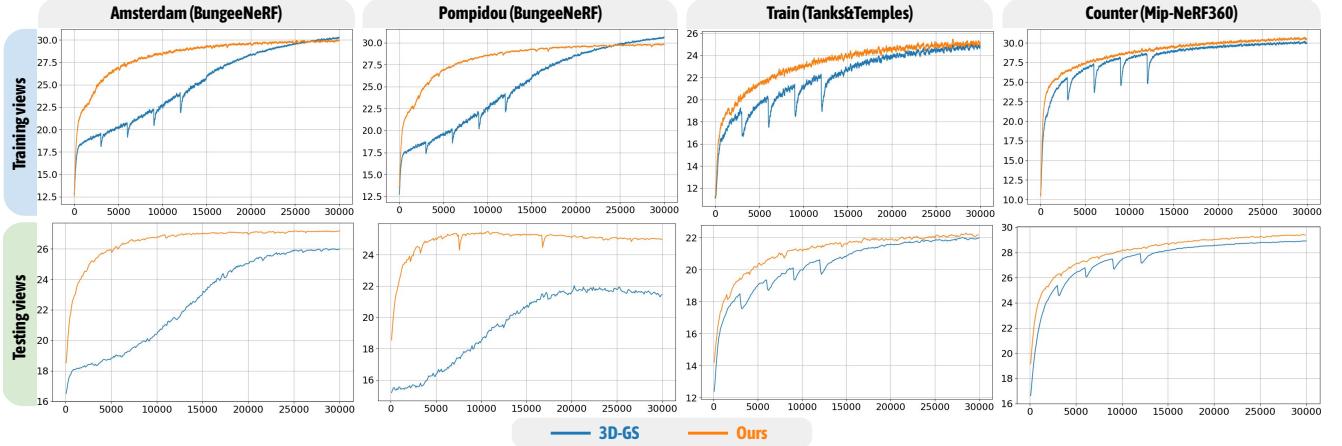
Figure 14. **PSNR curve of *Scaffold*-GS and 3D-GS [22] across diverse datasets [4, 17, 49].** We illustrate the variations in PSNR during the training process for both training and testing views. The orange curve represents *Scaffold*-GS, while the blue curve corresponds to 3D-GS. Our method not only achieves rapid convergence but also exhibits superior performance, marked by a significant rise in training PSNR and consistently higher testing PSNR, in contrast to 3D-GS.

Table 11. **PSNR scores for Tanks&Temples [23] and Deep Blending [18] scenes.**

| Method | Scenes | Truck | Train | Dr Johnson | Playroom |
|---|---|---|---|---|---|
| **3D-GS [22]** | | 25.19 | 21.10 | 28.77 | 30.04 |
| **Mip-NeRF360 [4]** | | 24.91 | 19.52 | 29.14 | 29.66 |
| **iNPG [31]** | | 23.26 | 20.17 | 27.75 | 19.48 |
| **Plenoxels [13]** | | 23.22 | 18.93 | 23.14 | 22.98 |
| **Ours** | | **25.77** | **22.15** | **29.80** | **30.62** |

Table 12. **LPIPS scores for Tanks&Temples [23] and Deep Blending [18] scenes.**

| Method | Scenes | Truck | Train | Dr Johnson | Playroom |
|---|---|---|---|---|---|
| **3D-GS [22]** | | 0.148 | 0.218 | 0.244 | **0.241** |
| **Mip-NeRF360 [4]** | | 0.159 | 0.354 | **0.237** | 0.252 |
| **iNPG [31]** | | 0.274 | 0.386 | 0.381 | 0.465 |
| **Plenoxels [13]** | | 0.335 | 0.422 | 0.521 | 0.499 |
| **Ours** | | **0.147** | **0.206** | 0.250 | 0.258 |

Table 13. **Storage size (MB) for Tanks&Temples [23] and Deep Blending [18] scenes.**

| Method | Scenes | Truck | Train | Dr Johnson | Playroom |
|---|---|---|---|---|---|
| **3D-GS [22]** | | 578 | 240 | 715 | 515 |
| **Ours** | | **107** | **66** | **69** | **63** |

Table 14. **PSNR scores for Synthetic Blender [30] scenes.**

| Method | Scenes | Mic | Chair | Ship | Materials | Lego | Drums | Ficus | Hotdog |
|---|---|---|---|---|---|---|---|---|---|
| **3D-GS [22]** | | 35.36 | **35.83** | 30.80 | 30.00 | **35.78** | 26.15 | 34.87 | 37.72 |
| **Ours** | | **37.25** | 35.28 | **31.17** | **30.65** | 35.69 | **26.44** | **35.21** | **37.73** |

Table 15. **Storage size (MB) for Synthetic Blender [30] scenes.**

| Method | Scenes | Mic | Chair | Ship | Materials | Lego | Drums | Ficus | Hotdog |
|---|---|---|---|---|---|---|---|---|---|
| **3D-GS [22]** | | 50 | 116 | 63 | 35 | 78 | 93 | 59 | 44 |
| **Ours** | | **12** | **13** | **16** | **18** | **13** | **35** | **11** | **8** |

Table 16. **PSNR scores for BungeeNeRF [49] and VR-NeRF [51] scenes.**

| Method | Scenes | Amsterdam | Bilbao | Pompidou | Quebec | Rome | Hollywood | Apartment | Kitchen |
|---|---|---|---|---|---|---|---|---|---|
| **3D-GS [22]** | | 25.74 | 26.35 | 21.20 | 28.79 | 23.54 | 23.25 | 28.48 | 29.40 |
| **Ours** | | **27.10** | **27.66** | **25.34** | **30.51** | **26.50** | **24.97** | **28.87** | **29.61** |

Table 17. **Storage size (MB) for BungeeNeRF [49] and VR-NeRF [51] scenes.**

| Method | Scenes | Amsterdam | Bilbao | Pompidou | Quebec | Rome | Hollywood | Apartment | Kitchen |
|---|---|---|---|---|---|---|---|---|---|
| **3D-GS [22]** | | 1453 | 1337 | 2129 | 1438 | 1626 | 1642 | 202 | 323 |
| **Ours** | | **243** | **197** | **230** | **166** | **200** | **182** | **48** | **90** |

the Amsterdam and Pompidou scenes in BungeeNeRF, we trained them with images at *three coarser scales* and evaluated them at a *novel finer scale*. The fact that 3D-GS achieved higher training PSNR but lower testing PSNR indicates its tendency to overfit at training scales.

convergence, enhanced robustness, and better generalization compared to 3D-GS, as evidenced by the rapid increase in training PSNR and higher testing PSNR. Specifically, for