

---

# SA-GS: Semantic-Aware Gaussian Splatting for Large Scene Reconstruction with Geometry Constraint

---

Anonymous Author(s)

Affiliation

Address

email

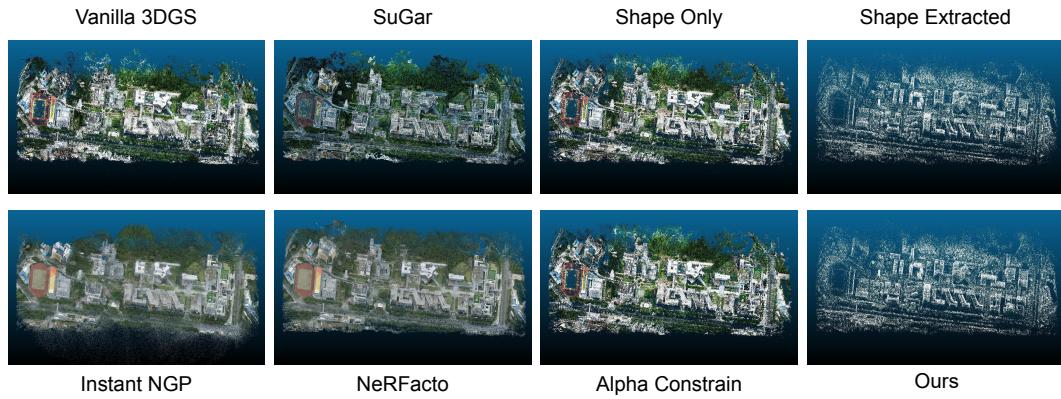


Figure 1: Qualitative comparison between our method and other 3DGS based methods. We proposed Shape constrain, alpha constrain and point cloud extraction in the current study. Quantitative ablation is shown in the right handside of the figure.

## Abstract

With the emergence of Gaussian Splats, recent efforts have focused on large-scale scene geometric reconstruction. However, most of these efforts either concentrate on memory reduction or spatial space division, neglecting information in the semantic space. In this paper, we propose a novel method, named SA-GS, for fine-grained 3D geometry reconstruction using semantic-aware 3D Gaussian Splats. Specifically, we leverage prior information stored in large vision models such as SAM and DINO to generate semantic masks. We then introduce a geometric complexity measurement function to serve as soft regularization, guiding the shape of each Gaussian Splat within specific semantic areas. Additionally, we present a method that estimates the expected number of Gaussian Splats in different semantic areas, effectively providing a lower bound for Gaussian Splats in these areas. Subsequently, we extract the point cloud using a novel probability density-based extraction method, transforming Gaussian Splats into a point cloud crucial for downstream tasks. Our method also offers the potential for detailed semantic inquiries while maintaining high image-based reconstruction results. We provide extensive experiments on publicly available large-scale scene reconstruction datasets with highly accurate point clouds as ground truth and our novel dataset. Our results demonstrate the superiority of our method over current state-of-the-art Gaussian Splats reconstruction methods by a significant margin in terms of geometric-based measurement metrics. Code and additional results are available on our anonymous [project page](#)

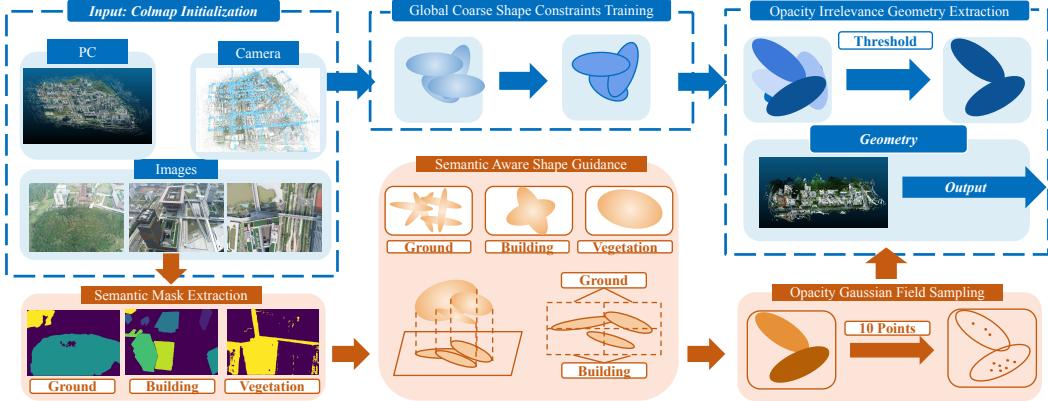


Figure 2: Overview: The blue section of the figure illustrates common methods for reconstructing geometrically aligned Gaussian Splats. The input for all Gaussian Splatting methods includes a COLMAP initialization consisting of images, camera positions, and SfM sparse point clouds. The output will be a traditional representation such as a mesh or point cloud, as shown in the right blue box. During training, in addition to the common image rendering loss, most methods encourage all 3D Gaussians to form a disk-like shape, as seen in [5] and [6]. After several training iterations, or at the end of the training process, other methods select a hard threshold for the alpha value and use the remaining Gaussians for geometric reconstruction. However, these hard constraints often result in poorer reconstruction, as demonstrated in our experiments. Instead of encouraging all Gaussians to adopt the same shape, our method uses semantic information to control the shape in detail. We first produce semantic masks for each input image, then extract shape information for each semantic group, and use this information to locally control the shape of each Gaussian. Additionally, we provide an opacity field sampling method that can dynamically allocate the desired number of points and ignore defective reconstruction parts.

## 22 1 Introduction

23 3D reconstruction is a transformative technology that converts real-world scenes into digital three-  
 24 dimensional models. As this technology often requires the transformation of multiple 2D images into  
 25 3D models, it finds numerous applications in urban planning, virtual reality (VR), and augmented  
 26 reality (AR). Various techniques have been employed to enhance the accuracy and efficiency of 3D  
 27 reconstruction, such as Neural Field-based methods [21][14][1][4][15][2] and Gaussian Splatting-  
 28 based methods [7][19][3][13][5].

29 Neural Rendering techniques, such as NeRF, often face challenges due to their lengthy training times  
 30 and the need for dense camera poses. These factors make them difficult to train, render and edit. In  
 31 contrast, 3D Gaussian Splatting (3DGS) [7] combines rasterization with novel view synthesis which  
 32 features rapid training and rendering speeds and exhibits high tolerance to sparse camera positions  
 33 and orientations.

34 Despite the merits of 3D Gaussian Splatting (3DGS) mentioned above, it often suffers from unrealis-  
 35 tic geometric reconstruction. Newly developed techniques primarily utilize both Signed Distance  
 36 Functions (SDF) and other mesh-based reconstruction methods in conjunction with Gaussian repre-  
 37 sentation to interactively train the system, ensuring that Gaussian splats adhere to the mesh surface.  
 38 This approach has demonstrated improved quantitative performance in geometric reconstruction.  
 39 However, this technique may generate random meshes when dealing with unbounded scenes, re-  
 40 sulting in a larger chamfer distance compared to the ground truth LiDAR point cloud. Additionally,  
 41 the works are based on the assumption that Gaussian splats will automatically adopt a disk-like  
 42 shape. By aligning the normal of the disk, the extracted mesh is expected to be smoother. The work  
 43 pushes this constraint further, as shown in the figure. These methods directly constrain the shape of  
 44 Gaussians, degenerating 3D Gaussian Splats to 2D or using a shape constraint to encourage a disk  
 45 shape Gaussian, facilitating mesh extraction by obtaining the normals of a 2D disk. However, curved  
 46 surfaces cannot be effectively represented using disks. The strategy to generate curved surfaces  
 47 typically results in the generation of a large number of 2D Gaussian Splats, leading to extremely high

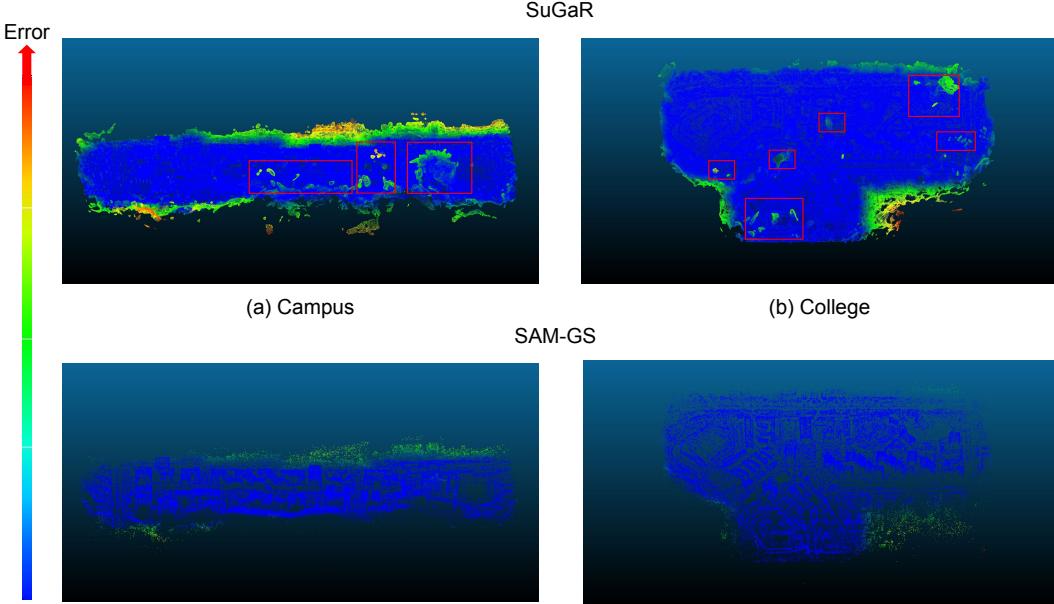


Figure 3: Explanation of fantasy-surface problem. In the first row of this figure, we display the results of using SuGaR [5] to reconstruct the Campus and College scenes from GauUsceneV2 [20]. Many surfaces incorrectly model the lighting conditions due to complex effects, such as how glass reflects sunlight at different angles and how clouds block sunlight. These imaginary surfaces that do not represent the true surface are regarded as fantasy surfaces. Our method, shown in the bottom rows, largely alleviates this problem, as evident in the figure. Another major source of geometric error occurs at the edges of unbounded scenes. However, this issue is common to all methods due to the sparsity of images at the edges and is not the focus of our current work.

48 memory consumption. Most geometric methods naively encourage every Gaussian to have relatively  
 49 high opacity to make Gaussians align with real objects’ surfaces. Although this naive approach can  
 50 work when there is stable lighting, complex lighting conditions can generate many small pieces of  
 51 meshes, which are even worse than vanilla Gaussian Splatting as shown in the figure. In the current  
 52 work, we regard this problem as a fantasy-surface problem. We further argue that simply encouraging  
 53 a disk-shaped and large alpha value Gaussian Splat for every Gaussian could detrimentally impact  
 54 overall reconstruction performance, as curved surfaces cannot be accurately represented by a disk-like  
 55 shape and lead to a fantasy-surface problem.

56 To solve the problem mentioned above, we introduce a semantic-aware masked model as shown in  
 57 the figure. We still take the same COLMAP input. But we control the shape of different Gaussians in  
 58 a fine-grained style according to their semantic attributes. We first obtain the semantic attributes using  
 59 GroundingSAM, and measure their geometric properties. We refer to all objects in the same category  
 60 as a semantic group. The geometric property of a semantic group is termed geometric perplexity. The  
 61 conclusion is that geometry complexity is related to the density of effective edges within a semantic  
 62 group. To reach this conclusion, we provide a detailed frequency domain analysis in the method  
 63 section.

64 A naive way to train is using the expected shape for each semantic group as a shape constraint to  
 65 detail control the shape of each Gaussian during training and merge all semantic groups into one  
 66 scene. However, this naive setting is not working due to the inconsistency of GroundingSAM. To  
 67 be specific, in some images, some semantic parts are regarded as ground while for the consecutive  
 68 image, the model will regard the same semantic part as something else as shown in the figure. We  
 69 then provide a robust loss function served as a soft regularization to encourage each Gaussian Splat  
 70 to form the shape we want. By first extending the geometric complexity idea to the whole model, we  
 71 can then calculate the lower bounds of the number of 3DGS we need to fully construct the scene.  
 72 Therefore, we provide a training strategy that iteratively decreases 3DGS during training, which  
 73 reduces the training memory consumption.



Figure 4: Explanation of Inconsistency problem. The semantic segmentation results are sometimes inconsistent with previous judgments. As shown in Figures (a) and (b), two tunnels are regarded as ground using GroundingSAM. However, in the images captured from a camera position immediately adjacent to them (Figures (c) and (d)), the left tunnel is not regarded as ground. This inconsistency between consecutive images is the primary cause of failure in naive reconstruction methods.

74 We have observed that geometric reconstruction often conflicts with lighting effects. Specifically,  
 75 when the reconstructed geometry aligns well with the point cloud, the image-based measurement  
 76 matrix such as SSIM, LPIPS, and PSNR, tends to perform poorly. Simply enforcing all Gaussian  
 77 has an alpha value larger than some specific opacity will lead to a fantasy-surface problem as we  
 78 mentioned before. To address this, we introduce the hierarchical probability density sampling strategy.  
 79 To be specific, low alpha value usually generates for complex lighting effect, by using or sampling  
 80 strategy, the extracted geometry will not sample the low alpha value area. Therefore the extracted  
 81 geometry is much better without harming the rendering result. Experiment results show that our  
 82 method surpasses the current SOTA method such as SuGaR and 2D Gaussian Splats by a large  
 83 margin.

## 84 2 Related Work

85 In this section, we introduce the Gaussian Splatting subsequent developments, particularly focusing  
 86 on large-scale, geometric reconstruction and object level semantic aware 3DGS.

### 87 2.1 Gaussian Splats on Large Scale Scene Reconstruction

88 Large-scale scene reconstruction faces challenges such as high memory usage, variable light-  
 89 ing, and sparse data. Notable techniques include CityGaussian[11], VastGaussian[9], and  
 90 HierarchicalGaussian[8], which employ a divide-and-conquer strategy, though their methods differ.  
 91 For instance, CityGaussian and VastGaussian segment based on camera visibility, while Hierarchical-  
 92 Gaussian uses a grid division. Techniques like EfficientGaussian[10] focus on memory efficiency,  
 93 introducing policies like gradient-sum thresholding and reducing low-impact Gaussians. Despite  
 94 these advancements, geometric accuracy is still under-validated, lacking high-accuracy datasets.  
 95 Recent datasets like GauUscene[20] and UrbanScene[12] offer more reliable data for testing. [20]  
 96 highlights the mismatch between image-based metrics (SSIM, LPIPS, PSNR) and geometric metrics  
 97 like Chamfer Distance, emphasizing the need for focused geometric assessments in our work.

### 98 2.2 Gaussian Splats on Geometric Reconstruction

99 Since the inception of 3D Gaussian Splatting [7], it has become a popular 3D representation method.  
 100 However, aligning splats geometrically is challenging. Key advancements include SuGaR[5] and  
 101 ScaffoldGaussian[13], which improve surface alignment using specialized loss functions and neural  
 102 networks, respectively. Techniques like 2DGGS[6] simplify splats to 2D to ease mesh extraction. The  
 103 novel SAGS approach enriches point clouds to enhance structural details. While these methods focus  
 104 on shape regularization, our work integrates shape and semantic insights using GroundedSAM[17]  
 105 for enhanced geometric reconstruction.

### 106 2.3 Object Level Semantic Aware Gaussian Splats

107 Semantic-aware approaches like LangSplat[16] and LeGaussian[18] leverage large vision models to  
 108 understand object semantics through embeddings. While effective at the object level, our approach

109 extends this understanding to large-scale scenes, utilizing semantic and geometric data to refine  
110 Gaussian Splats, showing marked improvements over purely geometric methods.

### 111 3 SA-GS

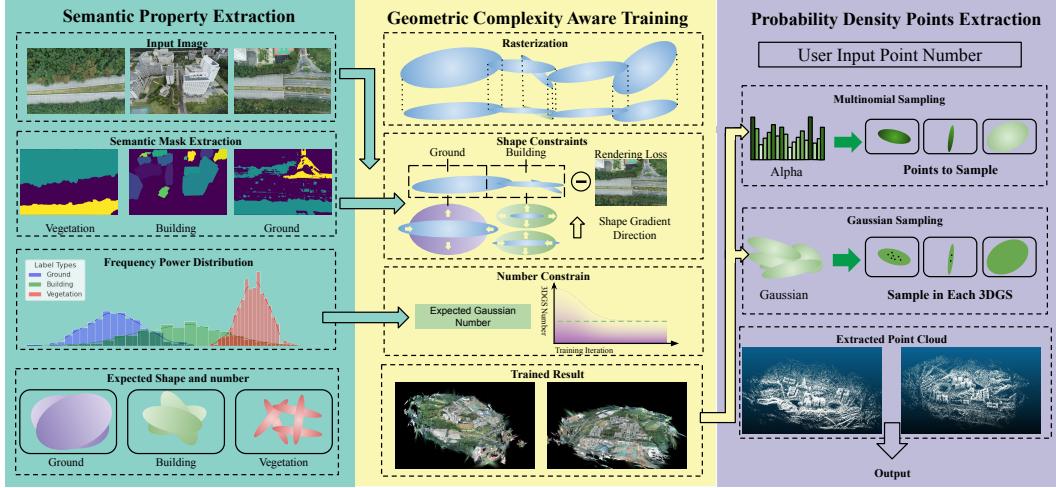


Figure 5: Method Overview: Our method pipeline consists of three main stages. Initially, we utilize the same input as vanilla Gaussian Splatting, but enhance it with semantic information extracted via Grounding SAM. Next, we assess the geometric complexity of each semantic group by calculating high-frequency power. Our geometric constraint is implemented through a soft regularization, facilitated by a semantic loss function. This guides the Gaussian shapes to match the expected shapes determined earlier. The rendering loss further refines the shape and attributes of the 3DGS, while the shape constraint, indicated by a negative sign, ensures alignment between rendered and real images. Controlling the shapes of different 3DGS is achieved by mapping their projected pixels onto the semantic map obtained earlier. Additionally, by reducing the number of low-opacity Gaussian splats to the expected count, we minimize GPU memory consumption during training. Finally, we offer a user-friendly point cloud extraction method via hierarchical probability density sampling. Initially, we create a multinomial distribution using the opacity values stored in each 3DGS. Then, based on user inputs and the multinomial distribution, we determine the number of points to sample from each Gaussian distribution. Detailed experimental results demonstrate significant improvements at each step, showcasing superior geometric reconstruction compared to current state-of-the-art methods.

112 Our method has three overall stages as shown in Fig.5. At the first stage, we transfer the input images  
113 to different semantic masked outputs utilizing GroundingSAM[17] and calculate the corresponding  
114 expected shape for each semantic group. At the second stage, we input the semantic map to start  
115 geometric complexity constraints training for Gaussian Splats. Notice that naively training different  
116 semantic groups separately will fail due to inconsistency, therefore we introduce a robust training  
117 and projection strategy at the training stage to solve this problem. In the third stage, we provide a  
118 hierarchical probability density extraction method for extracting the underlying geometry. In the  
119 following subsections, we will introduce each stage in detail.

#### 120 3.1 Geometric Property Extraction

121 The major goal of the current stage is to obtain the expected shape of 3DGS for each semantic group.  
122 In the end, we find that the expected shape of the Gaussian is related to the number of edges contained  
123 in a semantic group. The expected shape extraction can be divided into two main sections: one for  
124 semantic mask extraction and one for power distribution mapping.

##### 125 3.1.1 Semantic Mask Extraction.

126 We have  $N$  input images  $\mathcal{I} = \{I_1, I_2, \dots, I_N\}$ . We use GroundingSAM to get the masks of the  
127 images. The masked images are classified into different semantic groups. All masks are denoted as

128  $\mathcal{M} \in \mathbb{R}^{H \times W \times N}$ . We denote a pixel and position  $(x, y)$  in image  $i$  as  $M_i(x, y)$ . For every  $M_i(x, y)$ ,  
 129 it should be associated with a caption description or a default embedded value when GroundingSAM  
 130 cannot find the correspondence. Caption embedding is denoted as  $E(c)$ , where  $c$  is a text caption such  
 131 as "vegetation", "buildings", "road", and so on. During implementation, the embedding  $E(c)$  is set to  
 132 integer numbers for simplicity. We denote  $\forall M_i(x, y) = E(c_j)$  as a semantic group for caption  $c_j$ .  
 133 We have an overall  $C$  different captions. We then utilize the semantic group information as guidance  
 134 to constrain the shape of the Gaussian.

### 135 3.1.2 Power Distribution Mapping

136 Gaussian distributions are often regarded as low-pass filters due to their tendency to favor low  
 137 frequencies. In frequency domain analysis, the Gaussian Splatting model uses multiple low-frequency  
 138 models to represent a general spectrum in 3D space. To represent high-frequency information, such  
 139 as edges in 2D images or surfaces in 3D models, we use Gaussian splats with one relatively small  
 140 axis to achieve high frequency along that axis. For low-frequency regions, we use larger Gaussians.  
 141 Low-frequency areas cover a large space with fewer 3DGs, while high-frequency areas cover a small  
 142 space but require more 3DGs. For semantic groups with low high-frequency information, we use  
 143 disk-like shapes, and for groups with high-frequency signals, we use stick-like or dot-like Gaussians.  
 144 High-frequency signals, often represented by edges or corners in images, are modeled using Canny  
 145 edges. Mathematically, we define the expected ellipse with two aspect ratios.

$$146 a_1 = \frac{s_x}{s_z}, \quad a_2 = \frac{s_y}{s_z}, \quad (1)$$

146 where  $s_x, s_y, s_z$  are the scale along x, y, and z axis. The detailed mathematical derivation for how to  
 147 transfer from Gaussian shape to spectrum domain analysis and to edge extractor will be shown in the  
 148 technical appendix Sec.6.1. One can examine it in detail if one wants.

149 During the real implementation, we extract the edge energy utilizing Canny Edge instead of simple  
 150 high pass filter since Canny Edge can set the threshold for edge selection. Notice that Canny Edge is  
 151 not the only edge extractor can be applied in our method, basiclly any edge extractor can finish the  
 152 job. And we denote the edge count for image  $i$  that in semantic group  $j$  as  $e_{ij}$

153 We define the overall perplexity of a certain semantic group as this:

$$154 \mathbf{P}_j = \sum_i e_{ij}, \quad i \in \mathcal{M}. \quad (2)$$

154 This perplexity is for whole semantic group within a scene. Therefore, we can obtain the expected  
 155 Gaussian Splats number by multiply the edge number with a constant. This constant is usually  
 156 determined by the overlapping ratio of the input images. For unit perplexity, we need to divide by the  
 157 number of pixel that in one semantic group. We denote it as  $p_j$ .

158 Therefore, we have the following expectation for every Gaussian Splats in same semantic group:

$$159 \frac{1}{\mathbf{P}_i} = k_1 a_1, \quad \frac{1}{\mathbf{P}_i} = k_2 a_2, \quad k_1 > k_2. \quad (3)$$

### 159 3.2 Geometric Complexity Aware Training

160 To encourage Gaussian in different semantic group to align with the expected shape we have, we use  
 161 the following loss function as geometric complexity loss function.

$$162 \mathcal{L}_{gc} = \sigma(\frac{k_1}{\mathbf{P}_i} - a_1) + \sigma(\frac{k_2}{\mathbf{P}_i} - a_2). \quad (4)$$

162 As we mentioned before,we cannot directly train different semantic group separately due to the  
 163 inconsistency issue. Therefore, we modify the CUDA kernel to extract the project mean of each  
 164 Gaussian. We dynamically get the mask of Gaussian during training by attain the projected Gaussian  
 165 location on 2D masks. Then we assign the expected shape to each Gaussian online. In this way,  
 166 we circumvent the problem of inconsistency. In addition, since during rasterization process CUDA  
 167 keep the projected mean of each Gaussians, our modified method still keep the same training time as  
 168 memory consumption as vanilla Gaussian.

169 We adopt the same rendering loss as vanilla Nerf and use a hyper-parameter to tune the relationship  
 170 between those loss as following:

$$\lambda_{gc}\mathcal{L}_{gc} + \lambda_{dssim}\mathcal{L}_{dssim} + \lambda_{l1}\mathcal{L}_{l1}. \quad (5)$$

171 Where  $\mathcal{L}_{gc}$  is the geometry complexity loss.  $\mathcal{L}_{dssim}$  and  $\mathcal{L}_{l1}$  are the original rendering loss from  
 172 vanilla 3DGS. We set  $\lambda_{gc}, \lambda_{dssim}, \lambda_{l1}$  to 0.2, 0.2, and 0.6. We removes our Gaussian splats according  
 173 to ranking, and it gradually decrease to the expected number linearly utilizing after first 6000 iteration.

### 174 3.3 Probability Density Point Extraction

175 After obtained the trained Gaussian Spalts, we need to further extract the point cloud. Compare to  
 176 previous method such as GaussianPro[3], GauUscene[20] using mean of each Gaussian as extracted  
 177 point cloud we apply a hierarchical probability distribution sampling strategy that best fit the geometric  
 178 reconstruction requirement.

179 We posit that the probability of a position containing a point is proportional to the product of Gaussian  
 180 Density and opacity, given by:

$$\phi(x) = \sum_i \alpha_i \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)\right). \quad (6)$$

181 Notice the Equ.6 shows that the probability density function at position (x, y, z) is both proportional  
 182 to Gaussian Probability density function and alpha value. Therefore we construct a two stage  
 183 hierarchical sampling strategy as following. We first construct a multi-nominal distribution according  
 184 to alpha value:

$$P(i) = \frac{\alpha_i}{\sum \alpha}. \quad (7)$$

185 By utilizing multi nominal distribution, we can first sample as many points as we want according to  
 186 the alpha distribution. In reality, for index i with high alpha value, we might have multiple points  
 187 fall inside, while for other index with relatively low alpha value, there might be no point fall inside.  
 188 Then we sample the points within Gaussian. In this way we can sample as many points as we want.  
 189 Experiment results shows superiority in geometric reconstrucion when applying the above point  
 190 cloud extraction strategy.

## 191 4 Experiment

192 To validate our proposed method, we conducted intensive tests on both a public dataset, GauUScene  
 193 V2, covering over 6.5 square kilometers, and our own "Technology Campus" dataset, spanning 1.8  
 194 square kilometers. Additionally, we used the UrbanScene3D Polytech scene (1.7 square kilometers)  
 195 exclusively for image-based rendering validation due to its limited LiDAR data and lack of RGB  
 196 information. Collectively, these datasets cover more than 10 square kilometers, demonstrating our  
 197 method's effectiveness.

### 198 4.1 Implementation Detial and Metrics

199 We compared our method against several leading techniques, including SuGaR, Vanilla Gaussian  
 200 Splatting, InstantNGP, and NeRFacto, using the RTX 3090 for training. Our evaluation focused  
 201 on geometric reconstruction, primarily using Chamfer distance, but also included image-based  
 202 metrics such as PSNR, SSIM, and LPIPS for comprehensive assessment. The datasets and additional  
 203 validation results will be made available in the supplementary materials. We train our model use  
 204 RTX3090 and we set K1 to 3, k2 to 1.

### 205 4.2 Comparison

206 Since real-world LiDAR data is usually smaller than the reconstructed region using NeRF-based  
 207 and Gaussian-based methods, we crop our aligned reconstruction to the same size as the LiDAR  
 208 data to maintain fairness. Due to page limitations, we only show the qualitative comparison between

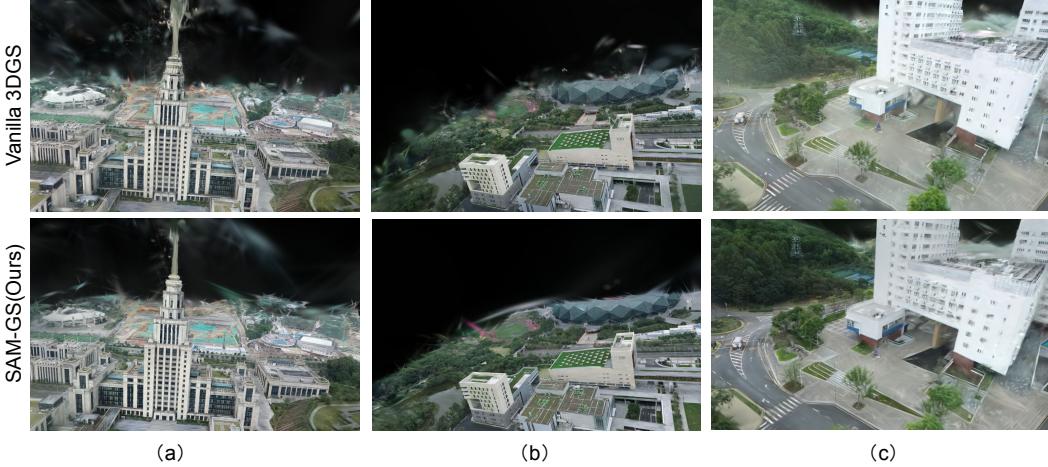


Figure 6: This the comparison between our method and vanilla Gaussian Splats. As one can see that from the figure shown above, our methods largely sharpening the edge of image. The tower shown in the figure (a) merges together and sharpened in our method, while in (b) figures, we eliminate the noise around the high building. While for the last group of pictures shows that our steadily alpha decreasing strategy is successful.

209 our method and the vanilla Gaussian. We display rendering results qualitatively in Fig. 6. We then  
 210 display our quantitative results on geometric alignment in Fig.1. Our geometric-based metric results  
 211 are shown in Tab.1. As for image-based metrics, we sample 10% of images from the whole dataset as  
 212 the testing dataset, with the remaining data used for training. Since our training procedure strictly  
 213 follows the GauScene V2 benchmark, all results on GauScene image-based metrics other than SA-GS  
 214 are directly borrowed from their report. We find using one card to train the NeRF model is much  
 215 faster than training it in parallel; therefore, we train all models using one RTX 3090. The detailed  
 216 comparison is shown in Tab.2 and Tab.3

Table 1: Comparison on Chamful distance and variance. This table displays the results obtained when testing two NeRF-based methods and 3DGS (3D Gaussian Splatting). For training and evaluation, we used SuGaR and the official Gaussian Splats code. Meanwhile, the NeRF Studio implementation was utilized for Instant-NGP and NeRFacto to conduct training and evaluation.

Method Scene   Metrics	Gaussian Splatting		SuGaR		SA-GS(Ours)		Instant NGP		NeRFacto	
	Mean ↓	Var ↓	Mean ↓	Var ↓	Mean ↓	Var ↓	Mean ↓	Var ↓	Mean ↓	Var ↓
Campus	0.044	0.117	0.081	0.163	<b>0.037</b>	<b>0.104</b>	0.083	0.265	0.054	0.169
Modern Building	0.040	<b>0.101</b>	0.066	0.140	0.040	0.108	0.065	0.186	<b>0.032</b>	0.140
Village	0.039	0.087	0.055	0.117	<b>0.035</b>	0.081	0.103	0.323	0.036	<b>0.077</b>
Residence	0.240	0.176	Nan	Nan	<b>0.203</b>	<b>0.148</b>	0.204	0.223	0.263	0.226
Russian Building	0.097	0.273	0.186	0.390	0.083	0.250	0.056	0.167	<b>0.041</b>	<b>0.155</b>
College	0.038	0.107	0.058	0.139	0.028	0.103	0.087	0.270	<b>0.024</b>	<b>0.082</b>
Technology Campus	0.073	0.229	0.146	0.339	<b>0.056</b>	<b>0.204</b>	0.057	0.166	0.093	0.270
Avg.	0.076	0.155	0.098	0.215	<b>0.068</b>	<b>0.143</b>	0.094	0.229	0.078	0.160

### 217 4.3 Ablation Study

218 In this section, we clearly display how each module in our model is used and its effect on the final  
 219 result. We have four groups of ablation studies. The first group involves shape constraint only,  
 220 followed by shape constraint with alpha constraint, which means steadily decreasing the number  
 221 of Gaussian Splats after the warm-up iteration. Notice that we did not apply hierarchical point  
 222 sampling for these two Gaussian experiments; instead, we directly use the mean of each Gaussian  
 223 to represent the extracted point cloud. The following experiments involve the extracted point cloud  
 224 from shape-only constraint and the extracted point cloud from Gaussian Splats with alpha constraint.  
 225 We use GauSceneV2 for testing. Since two of the representations are point clouds, we use Chamfer

Table 2: Image-based metrics result. The methods used are the same as mentioned above. We also provide additional training time for a detailed comparison. In the first six datasets, according to the GauUscene paper, they used four GPUs for training. Therefore, the last two datasets were trained on a single RTX 3090.

Method Scene   Metrics	Instant NGP				NeRFacto			
	PSNR ↑	SSIM ↑	LPIPS ↓	Time(GPU·min) ↓	PSNR ↑	SSIM ↑	LPIPS ↓	Time(GPU·min) ↓
Campus	20.76	0.516	0.817	220	17.70	0.455	0.779	1692
Modern Building	20.25	0.522	0.816	392	18.66	0.448	0.734	1704
Village	20.79	0.511	0.792	268	16.95	0.399	0.727	1788
Residence	18.64	0.453	0.856	348	15.05	0.364	0.879	1780
Russian Building	18.37	0.507	0.810	252	16.61	0.405	0.682	1716
College	19.64	0.551	0.820	276	17.28	0.462	0.781	1732
Technology Campus	19.58	0.510	0.720	<b>38</b>	17.78	0.463	0.805	685
polytech	17.83	0.494	0.843	<b>44</b>	15.55	0.460	0.928	751
Avg.	19.73	0.508	0.809	230	16.95	0.432	0.789	1481

PSNR ↑	Vanilla Gaussian			SuGaR			SA-GS(Ours)		
	SSIM ↑	LPIPS ↓	Time ↓	PSNR ↑	SSIM ↑	LPIPS ↓	Time ↓	PSNR ↑	SSIM ↑
24.76	<b>0.735</b>	<b>0.343</b>	<b>58</b>	23.02	0.601	0.506	104	<b>25.16</b>	0.730
<b>25.49</b>	<b>0.762</b>	<b>0.273</b>	64	22.51	0.572	0.497	108	25.21	0.739
<b>26.14</b>	<b>0.805</b>	<b>0.237</b>	<b>62</b>	22.78	0.619	0.461	98	25.64	0.783
<b>22.03</b>	<b>0.678</b>	<b>0.371</b>	<b>71</b>	20.97	0.533	0.607	119	21.26	0.629
<b>23.90</b>	<b>0.784</b>	<b>0.248</b>	<b>63</b>	21.58	0.618	0.450	103	23.86	0.771
<b>24.21</b>	<b>0.749</b>	<b>0.326</b>	<b>68</b>	22.02	0.588	0.514	123	24.14	0.724
<b>23.94</b>	<b>0.786</b>	<b>0.223</b>	69	21.47	0.556	0.431	146	23.05	0.741
22.31	0.772	0.273	77	20.98	0.569	0.487	173	<b>22.32</b>	<b>0.778</b>
<b>24.10</b>	<b>0.758</b>	<b>0.287</b>	<b>60</b>	21.92	0.528	0.494	121.8	23.82	0.736
								0.321	66.9

Table 3: This table shows the result we obtained using Gaussian Spalts based method

226 Distance as the measurement metric. The quantitative result is shown on 1. The quantitative result is  
227 shown in Tab.4

Table 4: Ablation results: As one can see, each part of our model is essential for achieving good geometric results. The result with alpha constraints and point cloud extraction yields the best performance in terms of CD Mean and CD variance.

Method Scene   Metrics	Gaussian Splatting		Shape Only		Alpha Constrain		Extracted Shpae Only		Extracted Alpha Constrain	
	Mean ↓	Var ↓	Mean ↓	Var ↓	Mean ↓	Var ↓	Mean ↓	Var ↓	Mean ↓	Var ↓
Campus	0.044	0.117	0.049	0.128	0.042	0.114	0.041	0.114	<b>0.037</b>	<b>0.104</b>
Modern Building	0.040	0.108	0.035	0.081	0.041	0.110	0.041	<b>0.108</b>	<b>0.040</b>	<b>0.108</b>
Village	0.039	0.087	0.043	0.097	0.038	0.088	0.038	0.087	<b>0.035</b>	<b>0.081</b>
Residence	0.240	0.176	0.265	0.188	0.231	0.168	0.209	0.156	<b>0.203</b>	<b>0.148</b>
Russian Building	0.097	0.273	0.105	0.296	0.092	0.267	0.093	0.278	<b>0.083</b>	<b>0.250</b>
College	0.038	0.107	0.047	0.122	0.034	0.105	0.032	0.108	<b>0.028</b>	<b>0.103</b>
Avg.	0.083	0.145	0.091	0.167	0.080	0.142	0.076	0.142	<b>0.071</b>	<b>0.132</b>

## 228 5 Conclusion

229 In our current work, we propose a semantic-aware geometric constraint algorithm that dynamically  
230 assigns expected shapes to Gaussian splats projected into different semantic groups. We present an  
231 algorithm capable of computing the geometric complexity of Gaussian splats based on spectrum  
232 analysis. Furthermore, we utilize geometric complexity measurement to determine the number of  
233 Gaussian splats. Subsequently, we introduce a hierarchical probability density sampling method  
234 that can extract as many points as desired by users while maintaining a dynamic alpha value to  
235 mitigate the fantasy surface problem. Additionally, we offer abundant experimental results. However,  
236 there are several drawbacks to our algorithm. Firstly, during training, we constrain the shape of all  
237 Gaussians that project onto the same pixel without explicitly ignoring Gaussians blocked by those  
238 with high opacity values before them. This may result in all Gaussians conforming to the shape of  
239 the semantic group that occupies the largest region in the scene. Secondly, our algorithm relies on  
240 key semantics provided by users, which may sometimes be absent. Thirdly, while the inconsistency  
241 between consecutive images can be addressed by our robust loss, the direct resolution of inconsistency  
242 in the 3D world itself has not been achieved.

243 **References**

- 244 [1] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-  
245 nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF  
246 Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022.
- 247 [2] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao  
248 Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In  
249 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14124–  
250 14133, 2021.
- 251 [3] Kai Cheng, Xiaoxiao Long, Kaizhi Yang, Yao Yao, Wei Yin, Yuexin Ma, Wenping Wang, and  
252 Xuejin Chen. Gaussianpro: 3d gaussian splatting with progressive propagation. *arXiv preprint  
253 arXiv:2402.14650*, 2024.
- 254 [4] Jonathan Granskog, Till N Schnabel, Fabrice Rousselle, and Jan Novák. Neural scene graph  
255 rendering. *ACM Transactions on Graphics (TOG)*, 40(4):1–11, 2021.
- 256 [5] Antoine Guédon and Vincent Lepetit. Sugar: Surface-aligned gaussian splatting for efficient 3d  
257 mesh reconstruction and high-quality mesh rendering. *arXiv preprint arXiv:2311.12775*, 2023.
- 258 [6] Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2d gaussian  
259 splatting for geometrically accurate radiance fields. *SIGGRAPH*, 2024.
- 260 [7] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian  
261 splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), July  
262 2023.
- 263 [8] Bernhard Kerbl, Andreas Meuleman, Georgios Kopanas, Michael Wimmer, Alexandre Lanvin,  
264 and George Drettakis. A hierarchical 3d gaussian representation for real-time rendering of very  
265 large datasets. *ACM Transactions on Graphics*, 43(4), July 2024.
- 266 [9] Jiaqi Lin, Zhihao Li, Xiao Tang, Jianzhuang Liu, Shiyong Liu, Jiayue Liu, Yangdi Lu, Xiaofei  
267 Wu, Songcen Xu, Youliang Yan, and Wenming Yang. Vastgaussian: Vast 3d gaussians for large  
268 scene reconstruction. In *CVPR*, 2024.
- 269 [10] Wenkai Liu, Tao Guan, Bin Zhu, Lili Ju, Zikai Song, Dan Li, Yuesong Wang, and Wei Yang.  
270 Efficientgs: Streamlining gaussian splatting for large-scale high-resolution scene representation.  
271 2024.
- 272 [11] Yang Liu, He Guan, Chuanchen Luo, Lue Fan, Junran Peng, and Zhaoxiang Zhang. City-  
273 gaussian: Real-time high-quality large-scale scene rendering with gaussians. *arXiv preprint  
274 arXiv:2404.01133*, 2024.
- 275 [12] Yilin Liu, Fuyou Xue, and Hui Huang. Urbanscene3d: A large scale urban scene dataset and  
276 simulator. 2021.
- 277 [13] Tao Lu, Mulin Yu, Lining Xu, Yuanbo Xiangli, Limin Wang, Dahua Lin, and Bo Dai. Scaffold-  
278 gs: Structured 3d gaussians for view-adaptive rendering. *arXiv preprint arXiv:2312.00109*,  
279 2023.
- 280 [14] Julien N. P. Martel, David B. Lindell, Connor Z. Lin, Eric R. Chan, Marco Monteiro, and  
281 Gordon Wetzstein. Acorn: Adaptive coordinate networks for neural scene representation. *ACM  
282 Trans. Graph. (SIGGRAPH)*, 40(4), 2021.
- 283 [15] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi,  
284 and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis.  
285 *Communications of the ACM*, 65(1):99–106, 2021.
- 286 [16] Minghan Qin, Wanhu Li, Jiawei Zhou, Haoqian Wang, and Hanspeter Pfister. Langsplat: 3d  
287 language gaussian splatting. *arXiv preprint arXiv:2312.16084*, 2023.
- 288 [17] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu  
289 Huang, Yukang Chen, Feng Yan, et al. Grounded sam: Assembling open-world models for  
290 diverse visual tasks. *arXiv preprint arXiv:2401.14159*, 2024.

- 291 [18] Jin-Chuan Shi, Miao Wang, Hao-Bin Duan, and Shao-Hua Guan. Language embedded 3d  
 292 gaussians for open-vocabulary scene understanding. *arXiv preprint arXiv:2311.18482*, 2023.
- 293 [19] Joanna Waczyńska, Piotr Borycki, Sławomir Tadeja, Jacek Tabor, and Przemysław Spurek.  
 294 Games: Mesh-based adapting and modification of gaussian splatting. *arXiv preprint  
 295 arXiv:2402.01459*, 2024.
- 296 [20] Butian Xiong, Nanjun Zheng, and Zhen Li. Gauu-scene v2: Expanse lidar image dataset  
 297 shows unreliable geometric reconstruction using gaussian splatting and nerf. *arXiv preprint  
 298 arXiv:2404.04880*, 2024.
- 299 [21] Xiaoshuai Zhang, Sai Bi, Kalyan Sunkavalli, Hao Su, and Zexiang Xu. Nerfusion: Fusing  
 300 radiance fields for large-scale scene reconstruction. In *Proceedings of the IEEE/CVF Conference  
 301 on Computer Vision and Pattern Recognition (CVPR)*, pages 5449–5458, June 2022.

## 302 6 Appendix

### 303 6.1 Mathematical Derivation of Spectrum, Edge, and Gaussian Shape relationship

304 Mathematically, we define the expected ellipse with two aspect ratios.

$$305 \quad a_1 = \frac{s_x}{s_z}, a_2 = \frac{s_y}{s_z}, \quad (8)$$

306 where  $s_x, s_y, s_z$  are the scale along x, y, and z axis. The detailed mathematical derivation for how to  
 307 transfer from Gaussian shape to spectrum domain analysis and to edge extractor will be shown in  
 308 the technical appendix. One can examine it in detail if one wants. For each masked image  $F(x, y)$   
 309 we transfer it to Fourier Domain as  $F(u, v)$ . A naive way to calculate how many Gaussian we need  
 310 to use is to find the mean magnitude of frequency. When high frequency region is large, the mean  
 311 magnitude of frequency will be large. To calculate magnitude of frequency we can first transfer  
 312  $F(u, v)$  in to polar coordinates  $f(\rho, \theta)$  and get the mean magnitude in Equ.9,  $f_\mu$  is magnitude of  
 313 frequency. By calculating the weighted mean along the frequency orientation, we might find the  
 314 geometric complexity.

$$315 \quad f_\mu = \int_0^{\frac{\pi}{2}} \int_0^\infty \rho \cdot |f(\rho, \theta)| d\rho d\theta \quad (9)$$

316 However, this naive magnitude statistics are not working well according to our experimental result.  
 317 The major contribution of  $f_\mu$  are actually low frequency information even though it has lower weight.  
 318 We further observe that high frequency signal higher than certain threshold shown in image usually  
 319 will lead to a dandified Gaussian Splats. Frequency higher than the threshold has the same effect  
 320 on Gaussian Shape. This observation leads to a high-pass like statistic strategy as shown in Equ.10.  
 321 Where  $\delta$  function is an impulse function.

$$322 \quad f_\mu = \int_0^{\frac{\pi}{2}} \int_0^\infty D(\rho) \cdot |f(\rho, \theta)| d\rho d\theta \quad D(\rho) = \delta(T - \rho) \quad (10)$$

323 Notice that the magnitude of spectrum  $|f(\rho, \theta)|$  is positively related to the energy of spectrum  
 324  $E(F(u, v)D(\rho))$ . Notice that the energy itself is always the same. As shown in the Equ.11, Where  
 325 the inverse Fourier transform of step function will be  $\delta$  function with a phase shift. In here we simply  
 326 use an impulse function to represent since we are calculating the real energy sum in the end. The  
 327 left-hand-side of equation simply illustrate the energy of an image passing through a high pass filter.

$$328 \quad E = \int_{-\infty}^{\infty} |f(x, y) * \delta(x, y)|^2 dx dy \quad F^{-1}(D) = \delta(\rho) \quad (11)$$

329 The effect of a high pass filter is actually an edge extraction kernel especially in image processing.  
 330 During the real implementation, we extract the edge energy utilizing Canny Edge instead of simple  
 331 high pass filter since Canny Edge can set the threshold for edge selection. And we denote the edge  
 332 count for image  $i$  that in semantic group  $j$  as  $e_{ij}$

329 We define the overall perplexity of a certain semantic group as this:

$$\mathbf{P}_j = \sum_i e_{ij}, \quad i \in \mathcal{M} \quad (12)$$

330 This perplexity is for whole semantic group within a scene. Therefore, we can obtain the expected  
331 Gaussian Splats number by multiply the edge number with a constant. This constant is usually  
332 determined by the overlapping ratio of the input images. That is when the overlapping ratio is large,  
333 the constant should be small. For unit perplexity, we need to divide by the number of pixel that in  
334 one semantic group. We denote it as  $p_j$ .

335 Therefore, we have the following expectation for every Gaussian Splats in same semantic group:

$$\frac{1}{\mathbf{p}_i} = k_1 a_1 \quad \frac{1}{P_i} = k_2 a_2 \quad k_1 > k_2 \quad (13)$$

336 **NeurIPS Paper Checklist**

337 **1. Claims**

338 Question: Do the main claims made in the abstract and introduction accurately reflect the  
339 paper's contributions and scope?

340 Answer: [Yes]

341 Justification: In abstraction, we state clearly what is semantic guidance is, and since no body  
342 use this inforamtion in reconstruction, we do it.

343 Guidelines:

- 344 • The answer NA means that the abstract and introduction do not include the claims  
345 made in the paper.
- 346 • The abstract and/or introduction should clearly state the claims made, including the  
347 contributions made in the paper and important assumptions and limitations. A No or  
348 NA answer to this question will not be perceived well by the reviewers.
- 349 • The claims made should match theoretical and experimental results, and reflect how  
350 much the results can be expected to generalize to other settings.
- 351 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
352 are not attained by the paper.

353 **2. Limitations**

354 Question: Does the paper discuss the limitations of the work performed by the authors?

355 Answer: [Yes]

356 Justification: In the conclusion, we state clearly the limitation. Such as alpha blocking  
357 problem, inconsitancy problem and so on. Our assumption is simple and common in nature  
358 image. That is the edge or high frequency information is rare.

359 Guidelines:

- 360 • The answer NA means that the paper has no limitation while the answer No means that  
361 the paper has limitations, but those are not discussed in the paper.
- 362 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 363 • The paper should point out any strong assumptions and how robust the results are to  
364 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
365 model well-specification, asymptotic approximations only holding locally). The authors  
366 should reflect on how these assumptions might be violated in practice and what the  
367 implications would be.
- 368 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
369 only tested on a few datasets or with a few runs. In general, empirical results often  
370 depend on implicit assumptions, which should be articulated.
- 371 • The authors should reflect on the factors that influence the performance of the approach.  
372 For example, a facial recognition algorithm may perform poorly when image resolution  
373 is low or images are taken in low lighting. Or a speech-to-text system might not be  
374 used reliably to provide closed captions for online lectures because it fails to handle  
375 technical jargon.
- 376 • The authors should discuss the computational efficiency of the proposed algorithms  
377 and how they scale with dataset size.
- 378 • If applicable, the authors should discuss possible limitations of their approach to  
379 address problems of privacy and fairness.
- 380 • While the authors might fear that complete honesty about limitations might be used by  
381 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
382 limitations that aren't acknowledged in the paper. The authors should use their best  
383 judgment and recognize that individual actions in favor of transparency play an impor-  
384 tant role in developing norms that preserve the integrity of the community. Reviewers  
385 will be specifically instructed to not penalize honesty concerning limitations.

386 **3. Theory Assumptions and Proofs**

387 Question: For each theoretical result, does the paper provide the full set of assumptions and  
388 a complete (and correct) proof?

389                  Answer: [NA]

390                  Justification: We although not derive a theory, but we have a mathematical derivation of how  
391                  to transfer geometric compexity, Gaussian Splats, and high frequency power in appendix.  
392                  And we proof it holds by abundant experiment.

393                  Guidelines:

- 394                  • The answer NA means that the paper does not include theoretical results.
- 395                  • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
396                  referenced.
- 397                  • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 398                  • The proofs can either appear in the main paper or the supplemental material, but if  
399                  they appear in the supplemental material, the authors are encouraged to provide a short  
400                  proof sketch to provide intuition.
- 401                  • Inversely, any informal proof provided in the core of the paper should be complemented  
402                  by formal proofs provided in appendix or supplemental material.
- 403                  • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 404                  4. Experimental Result Reproducibility

405                  Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
406                  perimental results of the paper to the extent that it affects the main claims and/or conclusions  
407                  of the paper (regardless of whether the code and data are provided or not)?

408                  Answer: [Yes]

409                  Justification: We state clearly what the hyper parameter is, and we provide our draft code in  
410                  the supplementary (there are many thing need to be done manually, hard to follow, but it  
411                  works). We also provide project page and source code publicly.

412                  Guidelines:

- 413                  • The answer NA means that the paper does not include experiments.
- 414                  • If the paper includes experiments, a No answer to this question will not be perceived  
415                  well by the reviewers: Making the paper reproducible is important, regardless of  
416                  whether the code and data are provided or not.
- 417                  • If the contribution is a dataset and/or model, the authors should describe the steps taken  
418                  to make their results reproducible or verifiable.
- 419                  • Depending on the contribution, reproducibility can be accomplished in various ways.  
420                  For example, if the contribution is a novel architecture, describing the architecture fully  
421                  might suffice, or if the contribution is a specific model and empirical evaluation, it may  
422                  be necessary to either make it possible for others to replicate the model with the same  
423                  dataset, or provide access to the model. In general, releasing code and data is often  
424                  one good way to accomplish this, but reproducibility can also be provided via detailed  
425                  instructions for how to replicate the results, access to a hosted model (e.g., in the case  
426                  of a large language model), releasing of a model checkpoint, or other means that are  
427                  appropriate to the research performed.
- 428                  • While NeurIPS does not require releasing code, the conference does require all submis-  
429                  sions to provide some reasonable avenue for reproducibility, which may depend on the  
430                  nature of the contribution. For example
  - 431                          (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
432                          to reproduce that algorithm.
  - 433                          (b) If the contribution is primarily a new model architecture, the paper should describe  
434                          the architecture clearly and fully.
  - 435                          (c) If the contribution is a new model (e.g., a large language model), then there should  
436                          either be a way to access this model for reproducing the results or a way to reproduce  
437                          the model (e.g., with an open-source dataset or instructions for how to construct  
438                          the dataset).
  - 439                          (d) We recognize that reproducibility may be tricky in some cases, in which case  
440                          authors are welcome to describe the particular way they provide for reproducibility.  
441                          In the case of closed-source models, it may be that access to the model is limited in  
442                          some way (e.g., to registered users), but it should be possible for other researchers  
443                          to have some path to reproducing or verifying the results.

444     **5. Open access to data and code**

445     Question: Does the paper provide open access to the data and code, with sufficient instruc-  
446     tions to faithfully reproduce the main experimental results, as described in supplemental  
447     material?

448     Answer: [Yes]

449     Justification: We provide code in the supplementary, and will continuously improve the code  
450     instruction soon. Our code and data is public available, as long as one do not use this dataset  
451     for commercial or military use.

452     Guidelines:

- 453       • The answer NA means that paper does not include experiments requiring code.
- 454       • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 455       • While we encourage the release of code and data, we understand that this might not be  
456       possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
457       including code, unless this is central to the contribution (e.g., for a new open-source  
458       benchmark).
- 459       • The instructions should contain the exact command and environment needed to run to  
460       reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 461       • The authors should provide instructions on data access and preparation, including how  
462       to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 463       • The authors should provide scripts to reproduce all experimental results for the new  
464       proposed method and baselines. If only a subset of experiments are reproducible, they  
465       should state which ones are omitted from the script and why.
- 466       • At submission time, to preserve anonymity, the authors should release anonymized  
467       versions (if applicable).
- 468       • Providing as much information as possible in supplemental material (appended to the  
469       paper) is recommended, but including URLs to data and code is permitted.

472     **6. Experimental Setting/Details**

473     Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
474     parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
475     results?

476     Answer: [Yes]

477     Justification: We provide code, and specify the hyperparameters clearly

478     Guidelines:

- 479       • The answer NA means that the paper does not include experiments.
- 480       • The experimental setting should be presented in the core of the paper to a level of detail  
481       that is necessary to appreciate the results and make sense of them.
- 482       • The full details can be provided either with the code, in appendix, or as supplemental  
483       material.

484     **7. Experiment Statistical Significance**

485     Question: Does the paper report error bars suitably and correctly defined or other appropriate  
486     information about the statistical significance of the experiments?

487     Answer: [Yes]

488     Justification: We do not have an error bar setting here, the scene reconstruction problem  
489     result mostly is fixed if one select the same picture and have the same colmap dataset.

490     Guidelines:

- 491       • The answer NA means that the paper does not include experiments.
- 492       • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
493       dence intervals, or statistical significance tests, at least for the experiments that support  
494       the main claims of the paper.

- 495           • The factors of variability that the error bars are capturing should be clearly stated (for  
 496           example, train/test split, initialization, random drawing of some parameter, or overall  
 497           run with given experimental conditions).  
 498           • The method for calculating the error bars should be explained (closed form formula,  
 499           call to a library function, bootstrap, etc.)  
 500           • The assumptions made should be given (e.g., Normally distributed errors).  
 501           • It should be clear whether the error bar is the standard deviation or the standard error  
 502           of the mean.  
 503           • It is OK to report 1-sigma error bars, but one should state it. The authors should  
 504           preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
 505           of Normality of errors is not verified.  
 506           • For asymmetric distributions, the authors should be careful not to show in tables or  
 507           figures symmetric error bars that would yield results that are out of range (e.g. negative  
 508           error rates).  
 509           • If error bars are reported in tables or plots, The authors should explain in the text how  
 510           they were calculated and reference the corresponding figures or tables in the text.

511           **8. Experiments Compute Resources**

512           Question: For each experiment, does the paper provide sufficient information on the com-  
 513           puter resources (type of compute workers, memory, time of execution) needed to reproduce  
 514           the experiments?

515           Answer: [Yes]

516           Justification: Theoretically, one RTX3090 will be enough. We do modify some part of  
 517           SuGaR preprocessing code to use less memory. The time is shown in mins

518           Guidelines:

- 519           • The answer NA means that the paper does not include experiments.  
 520           • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
 521           or cloud provider, including relevant memory and storage.  
 522           • The paper should provide the amount of compute required for each of the individual  
 523           experimental runs as well as estimate the total compute.  
 524           • The paper should disclose whether the full research project required more compute  
 525           than the experiments reported in the paper (e.g., preliminary or failed experiments that  
 526           didn't make it into the paper).

527           **9. Code Of Ethics**

528           Question: Does the research conducted in the paper conform, in every respect, with the  
 529           NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

530           Answer: [Yes]

531           Justification: [TODO]

532           Guidelines:

- 533           • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.  
 534           • If the authors answer No, they should explain the special circumstances that require a  
 535           deviation from the Code of Ethics.  
 536           • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
 537           eration due to laws or regulations in their jurisdiction).

538           **10. Broader Impacts**

539           Question: Does the paper discuss both potential positive societal impacts and negative  
 540           societal impacts of the work performed?

541           Answer: [Yes]

542           Justification: Done in introduction

543           Guidelines:

- 544           • The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 566 11. Safeguards

567 Question: Does the paper describe safeguards that have been put in place for responsible  
 568 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
 569 image generators, or scraped datasets)?

570 Answer: [Yes]

571 Justification: On our project page, we have license and disclaimer

572 Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 583 12. Licenses for existing assets

584 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
 585 the paper, properly credited and are the license and terms of use explicitly mentioned and  
 586 properly respected?

587 Answer: [Yes]

588 Justification: We cite properly

589 Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- 597           • If assets are released, the license, copyright information, and terms of use in the  
 598 package should be provided. For popular datasets, [paperswithcode.com/datasets](http://paperswithcode.com/datasets)  
 599 has curated licenses for some datasets. Their licensing guide can help determine the  
 600 license of a dataset.  
 601           • For existing datasets that are re-packaged, both the original license and the license of  
 602 the derived asset (if it has changed) should be provided.  
 603           • If this information is not available online, the authors are encouraged to reach out to  
 604 the asset's creators.

605 **13. New Assets**

606 Question: Are new assets introduced in the paper well documented and is the documentation  
 607 provided alongside the assets?

608 Answer: [Yes]

609 Justification: New assets is in our project page

610 Guidelines:

- 611           • The answer NA means that the paper does not release new assets.  
 612           • Researchers should communicate the details of the dataset/code/model as part of their  
 613 submissions via structured templates. This includes details about training, license,  
 614 limitations, etc.  
 615           • The paper should discuss whether and how consent was obtained from people whose  
 616 asset is used.  
 617           • At submission time, remember to anonymize your assets (if applicable). You can either  
 618 create an anonymized URL or include an anonymized zip file.

619 **14. Crowdsourcing and Research with Human Subjects**

620 Question: For crowdsourcing experiments and research with human subjects, does the paper  
 621 include the full text of instructions given to participants and screenshots, if applicable, as  
 622 well as details about compensation (if any)?

623 Answer: [NA]

624 Justification:

625 Guidelines:

- 626           • The answer NA means that the paper does not involve crowdsourcing nor research with  
 627 human subjects.  
 628           • Including this information in the supplemental material is fine, but if the main contribu-  
 629 tion of the paper involves human subjects, then as much detail as possible should be  
 630 included in the main paper.  
 631           • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
 632 or other labor should be paid at least the minimum wage in the country of the data  
 633 collector.

634 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
 635 Subjects**

636 Question: Does the paper describe potential risks incurred by study participants, whether  
 637 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
 638 approvals (or an equivalent approval/review based on the requirements of your country or  
 639 institution) were obtained?

640 Answer: [NA]

641 Justification:

642 Guidelines:

- 643           • The answer NA means that the paper does not involve crowdsourcing nor research with  
 644 human subjects.  
 645           • Depending on the country in which research is conducted, IRB approval (or equivalent)  
 646 may be required for any human subjects research. If you obtained IRB approval, you  
 647 should clearly state this in the paper.

- 648
- We recognize that the procedures for this may vary significantly between institutions
- 649 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
- 650 guidelines for their institution.
- 651
- For initial submissions, do not include any information that would break anonymity (if
- 652 applicable), such as the institution conducting the review.