UNIG: MODELLING UNITARY 3D GAUSSIANS FOR VIEW-CONSISTENT 3D RECONSTRUCTION

Anonymous authors

Paper under double-blind review

Abstract

In this work, we present UniG, a view-consistent 3D reconstruction and novel view synthesis model that generates a high-fidelity representation of 3D Gaussians from sparse images. Existing 3D Gaussians-based methods usually regress Gaussians per-pixel of each view, create 3D Gaussians per view separately, and merge them through point concatenation. Such a view-independent reconstruction approach often results in a view inconsistency issue, where the predicted positions of the same 3D point from different views may have discrepancies. To address this problem, we develop a DETR (DEtection TRansformer)-like framework, which treats 3D Gaussians as decoder queries and updates their parameters layer by layer by performing multi-view cross-attention (MVDFA) over multiple input images. In this way, multiple views naturally contribute to modeling a unitary representation of 3D Gaussians, thereby making 3D reconstruction more view-consistent. Moreover, as the number of 3D Gaussians used as decoder queries is irrespective of the number of input views, allow an arbitrary number of input images without causing memory explosion. Extensive experiments validate the advantages of our approach, showcasing superior performance over existing methods quantitatively (improving PSNR by 4.2 dB when trained on Objaverse and tested on the GSO benchmark) and qualitatively.

027 028 029

030

025

026

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

3D object reconstruction and novel view synthesis (NVS) are pivotal in computer vision and graphics, converting 2D images into detailed 3D structures in various applications such as robotics, augmented reality, virtual reality, medical imaging, archaeology, and more. Neural Radiance Fields
(NeRF) attempts (Xu et al., 2024; Mildenhall et al., 2020; Wang et al., 2021a; Chen et al., 2021; Yu
et al., 2021; Liu et al., 2024a; Xiong et al., 2024) are notable in 3D fields recently. However, their
progress is impeded by slow rendering speeds due to the implicit. Recently, as a semi-implicit representation, 3D Gaussian Splatting (3D GS) (Kerbl et al., 2023) has achieved remarkable optimization
speed and high-quality novel view rendering performance in representing objects or scenes.

However, many recent methods based on 3D GS techniques encounter the challenge of view incon-039 sistency. This issue arises due to imprecise depth estimations from each views leading to duplicated 040 representations of the same object regions within the 3D reconstructions from different perspectives. 041 For instance, MVGamba (Yi et al., 2024) treats images and 3D Gaussians as sequences in Mamba 042 (Gu & Dao, 2023; Dao & Gu, 2024) which leads to input view order induced view inconsistency. 043 Splatter Image (Szymanowicz et al., 2024) and LGM (Tang et al., 2024a) predict pixel-aligned 3D 044 Gaussians for each input view in the camera space of the corresponding input view. Then these 3D Gaussians are transformed from camera spaces of each view to the world space and naively merged 046 together to obtain the ultimate 3D Gaussians, as depicted in fig. 1(a). 047

However, such dimension lifting from 2D images to 3D Gaussians in different views are independent and lacks interactions among different views, thus it may result in a single object point being represented by multiple 3D Gaussians at different positions, leading to the aforementioned view-inconsistency issue (Yang et al., 2024; Dong & Wang, 2024).

To address this issue, we propose a **Uni**tary 3D Gaussians (UniG) representation. Inspired by Deformable DETR (Liu et al., 2023b; Li et al., 2024; Zhang et al., 2023; Liu et al., 2022; Li et al., 2023a) that treats the position and properties of bounding box (Bbox) as queries of the Transformer



Figure 1: (a) Previous methods such as LGM (Tang et al., 2024a) directly concatenate Gaussians from different views, leading to view inconsistency. (b) Our method employs a unitary set of 3D Gaussians, projecting them onto each view and integrating information across views for Gaussian updates. (c) Our approach significantly surpasses previous methods in the novel view synthesis task.

decoder, we develop a DETR-like Transformer encoder-decoder framework, which treats 3D Gaus-067 sians as decoder queries and updates their parameters layer by layer by performing cross-attention 068 over multiple input views as keys and values. To work over multi-view input, we propose a multi-069 view deformable attention (MVDFA) operation, where each 3D point fetches related information from multi-view 2D images simultaneously, effectively guaranteeing the consistency. More specif-071 ically, MVDFA utilize camera modulation techniques (Karras et al., 2019; Hong et al., 2024) to diversify queries based on views. The queries are linearly transformed to make difference in each 073 view, with the weights and bias trained from camera parameters. Such operation gives each view 074 its corresponding camera pose information. The view-specific queries are then used for performing 075 deformable attention over corresponding images. Although similar to DFA3D (Li et al., 2023a) and 076 BEVFormer (Li et al., 2022) in employing deformable attention in 3D with a point projection strat-077 egy, our model prioritizes multi-view distinctions (different qureies in different views) to achieve a more precise 3D representation. Further elaboration is available in section 2.

079 As the number of 3D Gaussians is usually very large, e.g. over 10,000, the self-attention operation 080 in a deformable Transformer decoder layer will demand a significant memory and computational 081 cost. To improve the efficiency, inspired by (Wang et al., 2021b), we introduce a 3D Spatial Effi-082 cient Self-Attention (SESA) approach, leveraging Fast Point Sampling (FPS) (Qi et al., 2017a) to 083 downsize the number of keys and values while preserving the number of queries. Moreover, directly regressing the positions of 3D Gaussians may lead to convergence challenges (see appendix A.4). 084 To address this problem, we utilize a coarse-to-fine framework, where a direct lift from 2D to 3D is 085 employed for every pixel in randomly selected input views at the coarse stage. Then, the 3D Gaussians from this stage serve as the initialization for the deformable Transformer-based refinement 087 network, facilitating meaningful projected positions and aiding in convergence. 088

- In summary, our contributions are as follows:
 - We propose UniG, a novel 3D object reconstruction and NVS algorithm which utilizes a unitary set of 3D Gaussians as queries in deformable Transformer. Such an approach allows all input views to contribute to the same 3D representation and effectively addresses the view inconsistency issue and supports arbitrary number of input views.
 - We propose to use MVDFA for tackling the multi-view fusion challenge, SESA for minimizing the memory usage in self-attention, and a coarse-to-fine framework for mitigating the convergence issue when directly regressing world coordinates of 3D Gaussians.
 - Both quantitative and qualitative experiments are conducted for evaluation. Our proposed method achieves the state-of-the-art performance on the commonly-used benchmark.
- 099 100 101

090

091

092

093

095

096

098

061

062

063

064

065 066

- 2 RELATED WORK
- 102

3D reconstruction from images Recently, various methods have been explored to reconstruct detailed 3D object from limited viewpoints. (Liu et al., 2024b;c; Tang et al., 2024b; Song et al., 2021a) view the problem as an image-conditioned generation task. Leveraging pretrained generative models like Rombach et al. (2022), they achieve realistic renderings of novel views. However, diffusion models require longer time to generate 3D with multi-step denoising process, thus limiting their applicability in real-time scenarios. Recent methodologies that rely on a single forward pro-

108 cess for 3D reconstruction, utilizing Neural Radiance Field (NeRF) (Mildenhall et al., 2020) as a 109 robust 3D representation, have demonstrated effective performance in the field of 3D reconstruction. 110 (Yu et al., 2021; Cao et al., 2022; Guo et al., 2022; Lin et al., 2022; Li et al., 2023b; Müller et al., 111 2022; Liu et al., 2024d; Wei et al., 2024; Tochilkin et al., 2024; Xu et al., 2024; Yu et al., 2021; 112 Wang et al., 2021a; Chen et al., 2021). However, due to the slow rendering speed of NeRF, it is being supplanted by a new, super-fast, semi-implicit representation-3D Gaussian Splatting (3D GS) 113 (Kerbl et al., 2023). Triplane-Gaussian (Zou et al., 2024), Gamba (Shen et al., 2024), and DIG3D 114 (Wu et al., 2024) make promising results on single image 3D reconstruction. Various techniques 115 such as SplatterImage (Szymanowicz et al., 2024), LGM (Tang et al., 2024a), pixelSplat (Charatan 116 et al., 2024), and MVSplat (Chen, Yuedong and Xu, Haofei and Zheng, Chuanxia and Zhuang, 117 Bohan and Pollefeys, Marc and Geiger, Andreas and Cham, Tat-Jen and Cai, Jianfei, 2024) have 118 extended the application of 3D Gaussian Splatting to multi-view scenarios. In these approaches, 119 each input view is processed to estimate 3D Gaussians specific to the view, followed by a simple 120 concatenation of the resulting 3D Gaussian assets from all views. GS-LRM (Zhang et al., 2024) and 121 GRM (grm, 2024) exhibit a model structure similar to LGM, resulting in notable accomplishments 122 through enhanced training processes and consequently more precise depth regression. Nevertheless, 123 these models adhere to the pipeline of predicting 3D Gaussians separately for each view, they demands substantial computational resources, particularly as the number of views grows, the number 124 of Gaussians scales linearly with the number of views. Furthermore, these methods are unable to 125 accommodate an arbitrary number of views as input. 126

127

Deformable Transformer in 3D DFA3D (Li et al., 2023a) and BEVFormer (Li et al., 2022) are 128 introduced to address the feature-lifting challenge in 3D detection and autonomous driving tasks. 129 They achieve notable performance enhancements by employing a deformable Transformer to bridge 130 the gap between 2D and 3D. DFA3D initially uses estimated depth to convert 2D feature maps to 3D, 131 sampling around reference points for deformable attention in each view. However, the 3D sampling 132 point design causes all projected 2D points to represent a singular point, neglecting view variations. 133 BEVFormer (Li et al., 2022) regards the Bird's-Eye-View (BEV) features as queries, projecting 134 the feature onto each input view. The Spatial Cross-Attention facilitates the fusion of BEV and 135 image spaces, though challenges persist sampling 4 height values per pillar in the BEV feature for 136 selecting 3D reference points may limit coverage, posing challenges in accurate keypoint selection 137 for the model. When contrasting DFA3D and BEVFormer with our MVDFA, a commonality lies 138 in projecting onto 3D regression targets to extract data from various image perspectives. However, 139 our model diverges by employing camera modulation to differentiate queries across views, enabling more specific information retrieval. 140

141 142

143 144

145

3 Methods

3.1 PRELIMENARIES OF 3D GS

146 3D GS (Kerbl et al., 2023) is a novel rendering method that can be viewed as an extension of point-147 based rendering methods (Kerbl et al., 2023; Chen & Wang, 2024). Hence, 3D Gaussians can serve 148 as effective 3D representations for efficient differentiable rendering. Each 3D Gaussian ellipsoid 149 can be described by $\mathbf{G} = \{ \mathrm{SH}, \boldsymbol{\mu}, \boldsymbol{\sigma}, \mathbf{R}, \mathbf{S} \}$. The color of 3D Gaussians is represented by spherical 150 harmonics (SH $\in \mathbb{R}^{12}$) while the geometry is described by the center positions $\mu \in \mathbb{R}^3$, shapes 151 (covariance matrix Σ), and opacity ($\sigma \in \mathbb{R}$) of ellipsoids (Zwicker et al., 2001; Kerbl et al., 2023). 152 Especially, the covariance matrix can be optimized through a combination of rotation and scaling for each ellipsoid as $\Sigma = \mathbf{RSS}^T \mathbf{R}^T$, where $\mathbf{R} \in \mathbb{R}^4$ (represented by quaternion) represents the 153 rotation and $\mathbf{S} \in \mathbb{R}^3$ contains the scales in three directions. 154

155 156

157

3.2 OUR METHOD

Overall framework As illustrated in fig. 2, our model follows an encoder-decoder framework in
 a coarse-to-fine manner. We employ unitary 3D Gaussian representation, which define a unitary set
 of 3D Gaussians in the world space no matter how many input views are given. During the coarse
 stage, one or more images are randomly selected as input for a simple encoder-only model to directly
 predict 3D Gaussians, supervised by a RGB loss. Subsequently, in the refinement network, all input



Figure 2: Overall Framework: In the coarse stage, 3D Gaussians are produced for each pixel of the sampled random views from the input data. In the refinement stage, 3D Gaussians from the coarse stage serves as the initialization for the refinement network. Multi-view features extracted by the feature extractor serves as keys and values of decoder. Queries are updated by the decoder layer with image features and the positions of the centers of 3D Gaussians. The final 3D Gaussian representation is regressed from the queries. MVDFA: multi-view deformable attention in section 3.2.2. SESA: spatial efficient self-attention in section 3.2.3. FFN: feed-forward network.

162

163

164

165 166

167

173 174

183

images undergo processing through an image encoder and a cross-view attention module to extract
 multi-view image features (section 3.2.1).

Each 3D Gaussian is then projected onto each view to query relevant features and update their respective parameters by query refinement decoder with multi-view deformable attention (MVDFA)(section 3.2.2). Spatially efficient self-attention is utilized to reduce computational and memory costs, enabling the utilization of more 3D Gaussians for object reconstruction (section 3.2.3). Moreover, the coarse-to-fine design aims to ensure that the initial positions of the center of 3D Gaussians are not too distant from the ground truth or outside the field of view to gurantee the training convergence (section 3.2.4). The training objective is detailed in section 3.2.5.

193 194

195

3.2.1 FEATURE EXTRACTOR

To extract image features from multi-view input, we utilize UNet (Ronneberger et al., 2015; Song 196 et al., 2021b), a widely employed feature extractor in 3D reconstruction tasks, as demonstrated in 197 Tang et al. (2024a); Szymanowicz et al. (2024). To enhance the network's understanding of the complete 3D object, multi-view cross-attention is employed to transfer information among views 199 right after the UNet block, activated when the number of input views exceeds one. In this config-200 uration, each input view acts as queries, while the concatenation (post-flattening) of the remaining 201 views serves as keys and values. To efficiently enable cross-attention across all views, we employ 202 shifted-window attention, as introduced in the Swin Transformer (Liu et al., 2021). This mechanism 203 reduces interactions by focusing on tokens within a local window, effectively reducing memory us-204 age for large input sequences. By processing tokens within a fixed window, shifted-window attention 205 effectively lowers the computational complexity, thereby enhancing the overall efficiency.

206 207

208

3.2.2 VIEW-AWARE QUERY REFINEMENT DECODER

Decoder structure In the decoder module, we employ a fixed number of queries $\mathbf{Q} \in \mathbb{R}^{N \times C}$ with N and C denote the number of Gaussians and the hidden dimension to model 3D Gaussians by associating queries with 3D Gaussian ellipsoid parameters G, including the center μ , opacity σ , rotation \mathbf{R} , scaling \mathbf{S} , and Spherical Harmonics \mathbf{SH} . As depicted in fig. 2, the queries navigate through multiple decoder layers, each including a multi-view deformable attention (MVDFA) (section 3.2.2) mechanism to leverage image features, a spatial efficient self-attention (SESA) (section 3.2.3) layer for inter-Gaussian interactions, and a feed-forward network (FFN). The functionality of a decoder layer can be summarized by eq. (1), where \mathbf{F} represents image features from different views and \mathbf{P}^l signifies reference points in the l-th layer.

 $\mathbf{Q}^{l+1} = \text{FFN}(\text{SESA}(\text{MVDFA}(\mathbf{Q}^l, \mathbf{P}^l, \mathbf{F}))$ (1)

Finally, queries are processed through a splatter head S to compute $\Delta \mathbf{G} = S(\mathbf{Q})$ for updating the 3D Gaussian parameters: $\mathbf{G}' = \mathbf{G} + \Delta \mathbf{G}$ (except for rotation, which is updated by multiplication). Here, all views contribute to unitary 3D Gaussians, emphasizing the most relevant features. This strategy effectively alleviates the view inconsistency issue and is computationally more efficient.

223

218

Multi-view deformable attention (MVDFA) The goal of MVDFA is to enhance the unified
 queries and Gaussian representations by integrating multi-view image features. Following origi nal design of DFA, trainable sampling points are employed on the image features to sample the most
 relevant image features as values (Carion et al., 2020; Zhu et al., 2021). The remaining problem
 lies in determining the sampling points and attention scores. Specifically, we project the 3D queries
 onto each image view and adjust it using camera modulation to account for view discrepancies. The
 sampling offsets and attention scores are subsequently obtained from the view-specific queries.

By leveraging the center μ of each 3D Gaussian from the previous layer, along with the corresponding camera poses π_i and intrinsic parameters K_i for the *i*-th image, we can compute UV coordinates \mathbf{P}_i by projecting the center coordinates of each 3D Gaussian onto the image plane of the *i*-th input image using the pinhole camera model (Forsyth & Ponce, 2003; Hartley & Zisserman, 2003): $\mathbf{P}_i = K_i \pi_i \mu$. In this context, both matrices K_i and π_i are expressed in homogeneous form. These UV coordinates in \mathbf{P}_i then function as the reference points for 2D deformable attention.

237 As depicted fig. 3 (a), in 3D, we have a set of queries associated with 3D Gaussian paremeters \mathbf{G} 238 while the queries for each image planes should to be adjusted to suit each view individually. To 239 tackle this issue, we start by using camera modulation with the adaptive layer norm (adaLN) (Hong et al., 2024; Karras et al., 2019; 2020; Viazovetskyi et al., 2020) to generate view-specific queries. 240 More information on this modulation is provided in fig. 3(b). Subsequently, a linear layer to predicts 241 the sampling offsets Δs for retriving images features as values and another linear layer to predicts 242 the attention scores α of the sampling points s. Following (Zhang et al., 2023; Zhu et al., 2021), 243 we compute attention scores directly from queries, omitting keys to streamline calculations. Then, 244 we apply the grid sampling algorithm with bilinear interpolation to extract image features at these 245 sampling points, which act as the values \mathbf{v} for cross attention. 246

Finally, for each input view, we compute the updated queries for each view using the attention scores α and sampled values v. The ultimate unitary queries are then computed as a weighted sum of individual view queries, with the weights calculated using an linear layer on the view-specific queries. Detailed pseudo code for our multi-view deformable cross-attention is available in fig. 3(b).

251 252

3.2.3 SPATIAL EFFICIENT SELF-ATTENTION (SESA)

253 Our multi-view deformable cross-attention mechanism demonstrates a superior efficiency in terms 254 of computational cost and memory usage. However, self-attention is computationally expensive, 255 especially with numerous 3D Gaussians. Updating each Gaussian with information from all oth-256 ers may not always be essential, as neighboring Gaussians often contain similar information. To tackle this problem, drawing inspiration from Wang et al. (2021b), we introduce a method to re-257 duce the number of keys and values while keeping the number of queries unchanged during self-258 attention. This selective update strategy enables each query to be updated with a subset of related 259 queries, effectively enhancing the information exchange efficiency. To ensure crucial information 260 flow, we leverage the Fast Point Sampling (FPS) algorithm from point cloud methodologies (Qi 261 et al., 2017a;b). By utilizing Gaussian centers μ to identify distant points for querying, we opti-262 mize memory usage while guaranteeing essential information sharing among Gaussians. Additional 263 details are in appendix A.1. 264

265 266

3.2.4 COARSE-TO-FINE MODEL

Locating Gaussian centers in the world space In Szymanowicz et al. (2024), the Gaussian centers are located in each input view's camera space, i.e. $\mu_{cam} = [x_{cam}, y_{cam}, z_{cam}] = [u_1 d + \Delta_x, u_2 d + \Delta_y, d + \Delta_z]$, where the center coordinates $x_{cam}, y_{cam}, z_{cam}$ are parameterized by the depth d and offset values $(\Delta_x, \Delta_y, \Delta_z)$. The depth d represents the length of a ray originating from the camera



283 284

287

288

289

290

291

292

(a) MVDFA on the n-th 3D Gaussian

(b) Pseudo code

Figure 3: MVDFA: \mathbf{Q}_n denotes the *n*-th unitary queries while \mathbf{q}_{ni} denotes the *n*-th query on the *i*-th view modulated by the *i*-th camera \mathbf{Cam}_i . Linear layers are used on \mathbf{q}_{ni} to compute the sampling offsets $\Delta \mathbf{s}_{ni}$ and attention score α_{ni} . The *n*-th 3D Gaussian is projected onto images, and surrounding sampling points $\mathbf{s}_{ni} = \mathbf{P}_{ni} + \Delta \mathbf{s}_{ni}$ are sampled using offsets $\Delta \mathbf{s}_{ni}$. Values \mathbf{v}_{ni} are image features sampled at \mathbf{s}_{ni} . The final query is calculated by the weighted sum of updated view-specific queries \mathbf{q}'_{ni} , where w_i is the weight calculated by a linear layer on \mathbf{q}'_{ni} . *B* is batch size, *I* is the number of views, *C* is the hidden dimension, *N* is the number of Gaussians, *pinhole_proj* is the projection from 3D to 2D with the pinhole model. **F** is the image feature with height *H* and weight *W*. *K* and π are camera intrinsics and extrinsics, respectively.

293 294

295 center. u_1, u_2 are the UV coordinates of the ray passing through the corresponding input image. 296 This design represents each point with multiple Gaussians, potentially introducing view inconsis-297 tency due to concatenation issues at various points caused by depth inaccuracies and tend to shortcut 298 input views (Wu et al., 2024). In our framework, we define unitary Gaussians in the world space, 299 project their centers to each input view for feature retrieval, as depicted in fig. 3(a). The centers 300 of Gaussians can be written as $\mu_{\text{world}} = [x_{\text{world}}, y_{\text{world}}, z_{\text{world}}]$. However, during the initial training phases, discrepancies between the 3D Gaussian centers and ground truth often result in imprecise 301 selection of image features at sampling points, presenting challenges for model convergence. 302

We employ a relative coordinate system, where the camera poses for all views are known. The initial input view is established as the world coordinates (with the camera pose represented by the identity matrix), and subsequently, all other views are transformed to align with these coordinates. This approach allows us to represent all 3D data within this consistent relative coordinate system.

307

Coarse-to-fine To address this issue, we utilize a coarse network that directly regress 3D Gaussian parameters with one or more randomly selected input images as input. The role of this network is to provide a coarse initialization of 3D Gaussians for the subsequent refinement network. We use the UNet architecture as the feature extractor to train the coarse network. Subsequently, we use this trained parameters to initialize the refinement stage and independently train the refinement network.

313 314

315

3.2.5 TRAINING OBJECTIVE

316 Building upon prior 3D Gaussian-based reconstruction approaches, we leverage the differentiable rendering implementation by Kerbl et al. (2023) to generate RGB images from the 3D Gaussians 317 produced by our model. For each object, we render 4 input views and 8 additional views (12 views in 318 total) for supervision. Furthermore, aligning with the methodologies ((Hong et al., 2024; Tang et al., 319 2024a)), we employ a RGB loss in eq. (2), which consists of both a mean square error loss \mathcal{L}_{MSE} and 320 a VGG-based LPIPS (Learned Perceptual Image Patch Similarity) loss (Zhang et al., 2018a) \mathcal{L}_{LPIPS} 321 to guide the rendered views. Here I_{pd} represents the rendered views supervised by the ground truth 322 images I_{at} . 323

$$\mathcal{L} = \mathcal{L}_{\text{MSE}}(I_{pd}, I_{gt}) + \lambda \mathcal{L}_{\text{LPIPS}}(I_{pd}, I_{gt})$$
(2)

³²⁴ 4 EXPERIMENTS

This section delves into experiment details and outcomes. In section 4.1, we give dataset specifics, evaluation metrics, and implementation details. section 4.2 delves into quantitative and qualitative results for sparse view novel view synthesis, along with the visualization of 3D Gaussian centers as a point cloud. section 4.3 provides a comparative analysis of processing speeds and memory costs. section 4.4 presents an ablation study. Lastly, The versatility of our model extends to tasks such as image-to-3D and text-to-3D generation using a diffusion model, detailed in section 4.5.

- 332
- 333 334

4.1 DATASET AND EXPERIMENT SETTINGS

Dataset We utilized a refined subset of the Objaverse LVIS dataset (Deitke et al., 2023) for training and validation. The training dataset comprised two sets of rendered images: one set featured 12
 random camera poses, while the other included input rendering images captured from fixed viewpoints (front, back, left, right). Supervision was provided from 32 random views spanning elevations between -30 to 30 degrees. The resolution of the rendered images was downscaled to 128 × 128.

To evaluate our model, we conducted tests on the Google Scanned Objects (GSO) benchmark. Two test sets were utilized: one with fixed-view inputs (e.g., front, left, back, right) at 0 degrees elevation, tested on 32 random views with elevations ranging from 0 to 30 degrees, and the other includes 25 random views with corresponding camera poses. Importantly, there are no constraints on the elevation of the rendered views. We refer to these test sets as GSO-random and GSO-fixed in our subsequent analysis. More details for dataset can be found in appendix A.2.

Evaluation metric We compute the peak signal-to-noise ratio (PSNR), structural similarity index
 (SSIM) (Wang et al., 2004), and perceptual distance (LPIPS) (Zhang et al., 2018b) between the ren dered images and the ground truth. Additionally, we offer visual representations of both the rendered
 images and the 3D Gaussian centers as a point cloud. More details are provided in appendix A.2.

351 352

353

354 355

356 357

358

359

360

361

362

364

365

366

367

368

369

370

371

372

373

346

4.2 COMPARISON WITH STATE-OF-THE-ART METHODS

We provide the comparison with the state-of-the-art methods in this section.

4.2.1 FIXED VIEW INPUT

Table 1: Quantitative results for inputting 4 views on GSOfixed dataset. *The results of MV-Gamba are cited from the paper as they do not provide code or a test set.

Splatter Image (Szymanowicz et al., 2024)	25.6241	0.9151	0.1517
LGM (Smar) (Tang et al., 2024a) LGM (Large) (Tang et al., 2024a)	26.2487	0.7829	0.2180
InstantMesh (Xu et al., 2024) MV-Gamba* (Yi et al., 2024)	23.0177 26.2500	0.8893 0.8810	0.0886

We evaluated recent multi-view reconstruction models using 4 views as input. Splatter Image (Szymanowicz et al., 2024) were trained with their native data loaders, adjusting inputs to 4 views and supervision to 12. LGM and InstantMesh were evaluated using the provided checkpoints, with "Small" indicating models tailored to 128 resolution and "Large" to 256 resolution. All models were assessed assuming the same number of training views.

table 1 showcases the performance of these methods in novel view synthesis using 4 fixed views (front, back, right, left) on the GSO-fixed dataset. Our model surpassed previous approaches in PSNR, SSIM, and LPIPS for novel view synthesis, with a significant improvement of approximately 4.2 dB in PSNR. Additional results for 6 and 8 view inputs are available in appendix A.3.2.

We present visualization results for novel view synthesis in fig. 4 and 3D Gaussian centers represented as point clouds in appendix A.3.1. In our experiments, with resolution of 128, the LGM model corresponds to the small version. Observations in the figures reveal view inconsistency in LGM and a lack of details in InstantMesh, whereas our model maintains both details and view consistency. Further visualizations at a resolution of 256 are accessible in appendix A.3.1.

382

383

384

385

386

387

388

389

390

391

392

393 394

395

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413 414

415

416

417

418 419

420 421

422

423

424

425

426

427

428

429

Table 2: Quantitative results for inputting 4 views on GSO-random dataset.

Our Model	26 3020	0.8323	0.1570
LGM	15.1113	0.8440	0.1592
Splatter Image	25.7660	0.8932	0.2575
Method	PSNR \uparrow	SSIM \uparrow	LPIPS

Previous methods (LGM and InstantMesh) usually rely on fixed views as input, as they align well with views generated from diffusion models like ImageDream (Wang & Shi, 2023). In real-world scenarios, users are more inclined to provide random views as input. table 2 displays the results when utilizing random 4 views as input on the GSO-random dataset. Notably, there is a performance drop observed in LGM and InstantMesh with random input views. appendix A.3.1 provides the visualization results. For Splatter Image, although the PSNR does

not reduced much, its SSIM and LPIPS reduced significantly. We provide more visualization in appendix A.3.1 fig. 14.

4.2.3 INFERENCE ON ARBITRARY NUMBER OF VIEWS



Figure 5: Quantitative results with random number of views as input. The model is trained with 4 random input views and tested with variate number of views.

Training costs for 3D methods are considerable, often requiring 32 NVIDIA A100 (80G) GPUs over multiple days. Additionally, memory costs for previous methods increase linearly with the number of views, presenting challenges for training models with varying input views. Therefore, a model supporting inference with any number of inputs while being trained on a fixed set, such as 4 views, would provide significant advantages.

Our model retains unitary 3D Gaussians in world coordinates, treating views as complementary sources without compromising overall 3D integrity. This enables adaptability to variable view counts during inference, despite training on a fixed number of views. fig. 5 showcases the results of training the model with 4 random views and testing it with different number of views. More views results are in appendix A.3.2 fig. 17.

While other methods demonstrate satisfactory performance with 4 views during inference, their effectiveness diminishes as the view count deviates from 4. In contrast, our model excels as the number of views increases. It is important to highlight that LGM is not part of this comparison due to its incapacity to handle variations in the number of views between the training and testing phases.

```
4.3 INFERENCE TIME AND MEMORY COST
```

Table 3: Inference time comparison. 3D: forward time, render: rendering time, inference: time of one forward and 32 rendering. Unit in seconds.

Method	$3D\downarrow$	Render \downarrow	Inference ↓
DreamGaussian	118.3245	0.0038	118.4461
InstantMesh	0.6049	0.6206	20.4641
LGM	1.6263	0.0090	1.9143
Our Model	0.6939	0.0019	0.7538

430 431 We performed inference time tests across different model types, including a diffusionbased method (DreamGaussian (Tang et al., 2024b)), a NeRF-based model (InstantMesh (Xu et al., 2024)), a previous Gaussian-based model (LGM (Tang et al., 2024a)), and our model, as detailed in table 3. Our model maintains a reduced number of Gaussians and achieves the fastest rendering speed.

In contrast to previous methods that compute 3D Gaussians per pixel per input view, our model retains a single 3D Gaussian irrespective



Figure 4: Novel views on GSO-fixed dataset for inputting 4 views with resolution 128.

of the number of views. While conventional methods exhibit linear memory expansion with additional views or higher image resolutions, our approach sustains a consistent memory overhead or experiences slight increments due to the marginally higher cost of the image feature extractor. This design theoretically enables our model to accommodate more input views and higher resolutions for enhanced outcomes, potentially circumventing the out-of-memory limitations encountered by other methods.

481 4.4 ABLATION STUDIES 482

472 473 474

475

476

477

478

479

480

table 4 illustrates an ablation study that evaluates different components of the model architecture. All
 the experiments are evaluated on the Objaverse validation dataset. Removing the coarse stage and
 initializing randomly (without any constraint) results in the lowest performance. This problem arises
 from utilizing image features around the projected 3D Gaussian center within each image view, po-

486 tentially causing zero features and steep gradients when projections extend beyond the image plane. 487 When the coarse stage is randomly initialized within the cone of vision (CoV), performance im-488 proves. To provide a more meaningful initialization, we incorporate a coarse stage to acquire the 489 approximate Gaussian locations, followed by a refinement stage. This refined initialization empow-490 ers our model to achieve superior performance. Moreover, removing cross-view attention leads to a moderate decrease in performance compared to the full model. Using only the coarse stage (UNet-491 based) slightly underperforms the full model. Furthermore, removing the camera modulation on 492 queries or use 3D sampling points instead of sampling on each view adversely impacts the results, 493 underscoring the critical significance of this view-specific design. The full model achieves the best 494 performance across all metrics, indicating that each component contributes positively to the overall 495 model effectiveness. Additionally, we offer details on hyperparameter selections in appendix A.4. 496

Method	$PSNR \uparrow$	SSIM \uparrow	LPIPS \downarrow
w/o coarse (ran. init.)	12.1213	0.6531	0.6224
w/o coarse (ran. init. in CoV)	22.6740	0.8711	0.2383
w/o cross view attention	25.3923	0.9013	0.1007
coarse stage only (UNet)	25.6033	0.9107	0.0930
w/o camera modulation	26.1328	0.9201	0.0883
3D sampling points	25.8392	0.9117	0.0945
Full model	26.5334	0.9344	0.0667

Moreover, we add the ablation study on the number of views or different views input in the coarse stage in appendix A.4 table 9. We also give more view inconsistency visualization problem by visualize center of Gaussians from each view in different colors, as shown in fig. 7. Furthermore, removing the background use masks for Splatter Image and LGM may slightly improve the perfor-

mance (fig. 16, table 6)

4.5 APPLICATIONS IN 3D GENERATION

Image-to-3D conversion represents a fundamental application in 3D generation. Following the methodology of LGM and InstantMesh (Tang et al., 2024a; Xu et al., 2024), we initially leverage a multi-view diffusion model, ImageDream (Wang & Shi, 2023), to generate four predetermined views. Subsequently, our model is utilized for 3D Gaussian reconstruction. A comparative analysis with LGM and InstantMesh is detailed in appendix A.3.2. We also showcase the quality results of our model on both the GSO dataset and in-the-wild images in appendix A.3.1.

Our model can also do the Text-to-3D task. To evaluate quality, we utilize MVDream (Shi et al., 2024) to generate a single image from a text prompt. Subsequently, a diffusion model is employed to produce multi-view images, which are then processed by our model to derive a 3D representation. A qualitative comparison of the text-to-3D generation is presented in appendix A.3.1.

523 524

525

497

498

499

500

501

504

505

506

507

508

509 510

511 512

513

514

515

516

517

5 CONCLUSION AND LIMITATION

In this paper, we have introduced a novel sparse view 3D reconstruction and novel view synthe-526 sis method. Initially, a fixed number of 3D Gaussians with predefined properties are initialized, 527 and each Gaussian ellipsoid is projected onto input image features extracted by a feature extractor. 528 We propose the MVDFA block to integrate image features surrounding the projected 3D Gaussians 529 from each view to refine the 3D Gaussians, employing a coarse-to-fine strategy to ensure robust 530 model convergence. Additionally, we develop a spatially efficient self-attention mechanism to min-531 imize computational costs, tackling view inconsistency and computational inefficiency. Our model 532 accommodates an arbitrary number of views as input and showcases its effectiveness through quan-533 titative and qualitative experiments compared to state-of-the-art methods trained on Objaverse and 534 tested on the GSO dataset. Furthermore, with the aid of an off-the-shelf diffusion model, our model 535 undertakes generation tasks such as image-to-3D and text-to-3D conversions. We present an ablation study elucidating the significance of each model component. While our model signifies a 536 notable advancement in sparse view 3D reconstruction, there are inherent limitations. Presently, 537 user-provided camera parameters, both camera poses and intrisics, are necessary for projecting 3D 538 Gaussians onto images, presenting potential challenges in 3D reconstruction. Addressing this issue stands as a focal point for future research.

540 REFERENCES

560

565

566

567

568

573

574

575

580

581

582

583

542	GRM: Large Gaussian Reconstruction Model for Efficient 3D Reconstruction and Generation, au-
543	thor=Xu Yinghao and Shi Zifan and Yifan Wang and Chen Hansheng and Yang Ceyuan and Peng
544	Sida and Shen Yujun and Wetzstein Gordon, 2024.

- Ang Cao, Chris Rockwell, and Justin Johnson. FWD: Real-time novel view synthesis with forward warping and depth. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 15713–15724, 2022.
- 548
 549
 550
 551
 Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In European conference on computer vision (ECCV), 2020.
- David Charatan, Sizhe Lester Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 19457–19467, 2024.
- Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. MVSNeRF: Fast Generalizable Radiance Field Reconstruction from Multi-View Stereo. In 2021 IEEE/CVF International Conference on Computer Vision (ICCV), Oct 2021. doi: 10. 1109/iccv48922.2021.01386. URL http://dx.doi.org/10.1109/iccv48922.2021.
 01386.
- Guikun Chen and Wenguan Wang. A Survey on 3D Gaussian Splatting. <u>ArXiv</u>, 2024.
- 562 Chen, Yuedong and Xu, Haofei and Zheng, Chuanxia and Zhuang, Bohan and Pollefeys, Marc and
 563 Geiger, Andreas and Cham, Tat-Jen and Cai, Jianfei. Mvsplat: Efficient 3d gaussian splatting
 564 from sparse multi-view images. European Conference on Computer Vision, 2024.
 - Tri Dao and Albert Gu. Transformers are SSMs: Generalized Models and Efficient Algorithms Through Structured State Space Duality. In International Conference on Machine Learning (ICML), 2024.
- Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig
 Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In <u>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern</u> Recognition (CVPR), pp. 13142–13153, 2023.
 - Jiahua Dong and Yu-Xiong Wang. Vica-nerf: View-consistency-aware 3d editing of neural radiance fields. Advances in Neural Information Processing Systems (NIPS), 36, 2024.
- Laura Downs, Anthony Francis, Nate Koenig, Brandon Kinman, Ryan Hickman, Krista Reymann, Thomas B McHugh, and Vincent Vanhoucke. Google scanned objects: A high-quality dataset of 3d scanned household items. In <u>2022 International Conference on Robotics and Automation</u> (ICRA), pp. 2553–2560. IEEE, 2022.
 - David A Forsyth and Jean Ponce. A Modern Approach. <u>Computer vision: a modern approach</u>, 17: 21–48, 2003.
 - Albert Gu and Tri Dao. Mamba: Linear-Time Sequence Modeling with Selective State Spaces. arXiv preprint arXiv:2312.00752, 2023.
- Pengsheng Guo, Miguel Angel Bautista, Alex Colburn, Liang Yang, Daniel Ulbricht, Joshua M.
 Susskind, and Qi Shan. Fast and Explicit Neural View Synthesis. In 2022 IEEE/CVF Winter
 Conference on Applications of Computer Vision (WACV), Jan 2022. doi: 10.1109/wacv51458.
 2022.00009. URL http://dx.doi.org/10.1109/wacv51458.2022.00009.
- Richard Hartley and Andrew Zisserman. <u>Multiple View Geometry in Computer Vision</u>. Cambridge university press, 2003.
- Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli,
 Trung Bui, and Hao Tan. LRM: Large Reconstruction Model for Single Image to 3D. In
 International Conference on Learning Representations (ICLR), 2024.

594 Tero Karras, Samuli Laine, and Timo Aila. A Style-Based Generator Architecture for Generative 595 Adversarial Networks. In Proceedings of the IEEE/CVF conference on computer vision and 596 pattern recognition (CVPR), pp. 4401-4410, 2019. 597 Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyz-598 ing and Improving the Image Quality of StyleGAN. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR), pp. 8110-8119, 2020. 600 601 Bernhard Kerbl, Georgios Kopanas, Thomas Leimkuehler, and George Drettakis. 3D Gaussian 602 Splatting for Real-Time Radiance Field Rendering. In ACM Transactions on Graphics (TOG), 603 2023. 604 605 Diederik Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In International Conference on Learning Representations (ICLR), San Diega, CA, USA, 2015. 606 607 Hongyang Li, Hao Zhang, Zhaoyang Zeng, Shilong Liu, Feng Li, Tianhe Ren, and Lei Zhang. 608 DFA3D: 3D Deformable Attention For 2D-to-3D Feature Lifting. In Proceedings of the 609 IEEE/CVF international conference on computer vision (ICCV), 2023a. 610 611 Hongyang Li, Hao Zhang, Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, and Lei Zhang. 612 TAPTR: Tracking Any Point with Transformers as Detection. In Proceedings of the IEEE/CVF European Conference on Computer Vision (ECCV), 2024. 613 614 Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan 615 Sunkavalli, Greg Shakhnarovich, and Sai Bi. Instant3D: Fast Text-to-3D with Sparse-View 616 Generation and Large Reconstruction Model. The International Conference on Learning 617 Representations (ICLR), 2023b. 618 619 Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng 620 Dai. BEVFormer: Learning Bird's-Eye-View Representation from Multi-Camera Images via 621 Spatiotemporal Transformers. European conference on computer vision (ECCV), 2022. 622 Kai-En Lin, Yen-Chen Lin, Wei-Sheng Lai, Tsung-Yi Lin, Yichang Shih, and Ravi Ramamoor-623 thi. Vision Transformer for NeRF-Based View Synthesis from a Single Input Image. In 2023 624 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), 2022. 625 626 Kenkun Liu, Derong Jin, Ailing Zeng, Xiaoguang Han, and Lei Zhang. A Comprehensive Bench-627 mark for Neural Human Radiance Fields. Advances in Neural Information Processing Systems 628 (NIPS), 36, 2024a. 629 Minghua Liu, Ruoxi Shi, Linghao Chen, Zhuoyang Zhang, Chao Xu, Xinyue Wei, Hansheng Chen, 630 Chong Zeng, Jiayuan Gu, and Hao Su. One-2-3-45++: Fast Single Image to 3D Objects with Con-631 sistent Multi-View Generation and 3D Diffusion. In Proceedings of the IEEE/CVF Conference 632 on Computer Vision and Pattern Recognition, pp. 10072–10083, 2024b. 633 634 Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Mukund Varma T, Zexiang Xu, and Hao Su. 635 One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. In 636 Conference on Neural Information Processing Systems (NIPS), 2024c. 637 Minghua Liu, Chong Zeng, Xinyue Wei, Ruoxi Shi, Linghao Chen, Chao Xu, Mengqi Zhang, 638 Zhaoning Wang, Xiaoshuai Zhang, Isabella Liu, Hongzhi Wu, and Hao Su. Meshformer: High-639 quality mesh generation with 3d-guided reconstruction model. Conference on Neural Information 640 Processing Systems (NIPS), 2024d. 641 642 Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. 643 Zero-1-to-3: Zero-shot One Image to 3D Object. In IEEE/CVF International Conference on 644 Computer Vision (ICCV), 2023a. 645 Shilong Liu, Feng Li, Hao Zhang, Xiao Yang, Xianbiao Qi, Hang Su, Jun Zhu, and Lei Zhang. 646 DAB-DETR: Dynamic anchor boxes are better queries for DETR. In International Conference on 647 Learning Representations (ICLR), 2022.

648 649 650	Siyi Liu, Tianhe Ren, Jia-Yu Chen, Zhaoyang Zeng, Hao Zhang, Feng Li, Hongyang Li, Jun Huang, Hang Su, Jun-Juan Zhu, and Lei Zhang. Stable-DINO: Detection Transformer with Stable Matching. In <u>IEEE/CVF International Conference on Computer Vision (ICCV)</u> , 2023b.
652 653 654	Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin Transformer: Hierarchical vision transformer using shifted windows. In <u>Proceedings of the</u> <u>IEEE/CVF international conference on computer vision (ICCV)</u> , pp. 10012–10022, 2021.
655 656 657	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. In <u>The</u> <u>European Conference on Computer Vision (ECCV)</u> , 2020.
658 659 660	Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant Neural Graphics Prim- itives with a Multiresolution Hash Encoding. In <u>ACM Transactions on Graphics (SIGGRAPH)</u> , 2022.
662 663 664 665	Sharan Narang, Gregory Diamos, Erich Elsen, Paulius Micikevicius, Jonah Alben, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. Mixed Precision Training. In <u>6th international conference on learning representations (ICLR)</u> , volume 1, pp. 14, 2018.
666 667 668	Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation. In <u>Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)</u> , pp. 652–660, 2017a.
669 670 671 672	Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space. <u>Advances in neural information processing systems (NIPS)</u> , 30, 2017b.
673 674 675	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High- resolution image synthesis with latent diffusion models. In <u>Proceedings of the IEEE/CVF</u> <u>conference on computer vision and pattern recognition (CVPR)</u> , pp. 10684–10695, 2022.
676 677 678	Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional Networks for Biomed- ical Image Segmentation. In <u>Medical image computing and computer-assisted intervention</u> (MICCAI), 2015.
679 680 681 682	Qiuhong Shen, Zike Wu, Xuanyu Yi, Pan Zhou, Hanwang Zhang, Shuicheng Yan, and Xinchao Wang. Gamba: Marry gaussian splatting with mamba for single view 3d reconstruction. <u>arXiv</u> preprint arXiv:2403.18795, 2024.
683 684 685	Yichun Shi, Peng Wang, Jianglong Ye, Long Mai, Kejie Li, and Xiao Yang. MVDream: Multi-view Diffusion for 3D Generation. <u>The International Conference on Learning Representations (ICLR)</u> , 2024.
686 687 688	Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising Diffusion Implicit Models. In International Conference on Learning Representations (ICLR), 2021a.
689 690 691 692	Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-Based Generative Modeling through Stochastic Differential Equations. <u>The International Conference on Learning Representations (ICLR)</u> , 2021b. URL https: //openreview.net/forum?id=PxTIG12RRHS.
693 694 695	Stanislaw Szymanowicz, Christian Rupprecht, and Andrea Vedaldi. Splatter Image: Ultra-Fast Single-View 3D Reconstruction. In <u>IEEE/CVF Conference on Computer Vision and Pattern</u> <u>Recognition (CVPR)</u> , 2024.
696 697 698 699	Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, Gang Zeng, and Ziwei Liu. LGM: Large Multi-View Gaussian Model for High-Resolution 3D Content Creation. In European Conference on Computer Vision, pp. 1–18. Springer, 2024a.
700 701	Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. DreamGaussian: Generative Gaussian Splatting for Efficient 3D Content Creation. In <u>International Conference on Learning</u> Representations (ICLR), 2024b.

702 Dmitry Tochilkin, David Pankratz, Zexiang Liu, Zixuan Huang, , Adam Letts, Yangguang Li, Ding 703 Liang, Christian Laforte, Varun Jampani, and Yan-Pei Cao. TripoSR: Fast 3D Object Reconstruc-704 tion from a Single Image. arXiv preprint arXiv:2403.02151, 2024. 705 Yuri Viazovetskyi, Vladimir Ivashkin, and Evgeny Kashin. StyleGAN2 Distillation for Feed-706 forward Image Manipulation. In Proceedings of the IEEE/CVF European Conference on 707 Computer Vision (ECCV), pp. 170-186. Springer, 2020. 708 709 Peng Wang and Yichun Shi. ImageDream: Image-Prompt Multi-view Diffusion for 3D Generation. 710 arXiv preprint arXiv:2312.02201, 2023. 711 Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul Srinivasan, Howard Zhou, Jonathan T. 712 Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. IBRNet: Learn-713 ing Multi-View Image-Based Rendering. In 2021 IEEE/CVF Conference on Computer Vision 714 and Pattern Recognition (CVPR), Jun 2021a. doi: 10.1109/cvpr46437.2021.00466. URL 715 http://dx.doi.org/10.1109/cvpr46437.2021.00466. 716 717 Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, 718 and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction with-719 out convolutions. In Proceedings of the IEEE/CVF international conference on computer vision 720 (ICCV), pp. 568-578, 2021b. 721 Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image Quality Assessment: 722 From Error Visibility to Structural Similarity. IEEE transactions on image processing, 13(4): 723 600-612, 2004. 724 725 Xinyue Wei, Kai Zhang, Sai Bi, Hao Tan, Fujun Luan, Valentin Deschaintre, Kalyan Sunkavalli, Hao Su, and Zexiang Xu. Meshlrm: Large reconstruction model for high-quality mesh. arXiv 726 preprint arXiv:2404.12385, 2024. 727 728 Jiamin Wu, Kenkun Liu, Han Gao, Xiaoke Jiang, and Lei Zhang. DIG3D: Marrying Gaussian 729 Splatting with Deformable Transformer for Single Image 3D Reconstruction. arXiv preprint 730 arXiv:2404.16323, 2024. 731 732 Zhangyang Xiong, Chenghong Li, Kenkun Liu, Hongjie Liao, Jianqiao Hu, Junyi Zhu, Shuliang 733 Ning, Lingteng Qiu, Chongjie Wang, Shijie Wang, et al. MVHumanNet: A Large-scale Dataset of Multi-view Daily Dressing Human Captures. In Proceedings of the IEEE/CVF Conference on 734 Computer Vision and Pattern Recognition (CVPR), pp. 19801–19811, 2024. 735 736 Jiale Xu, Weihao Cheng, Yiming Gao, Xintao Wang, Shenghua Gao, and Ying Shan. InstantMesh: 737 Efficient 3D Mesh Generation from a Single Image with Sparse-view Large Reconstruction Mod-738 els. arXiv preprint arXiv:2404.07191, 2024. 739 Jiayu Yang, Ziang Cheng, Yunfei Duan, Pan Ji, and Hongdong Li. ConsistNet: Enforcing 3D 740 Consistency for Multi-view Images Diffusion. In Proceedings of the IEEE/CVF Conference on 741 Computer Vision and Pattern Recognition (CVPR), pp. 7079-7088, 2024. 742 743 Xuanyu Yi, Zike Wu, Qiuhong Shen, Qingshan Xu, Pan Zhou, Joo-Hwee Lim, Shuicheng Yan, 744 Xinchao Wang, and Hanwang Zhang. MVGamba: Unify 3D Content Generation as State Space 745 Sequence Modeling. arXiv preprint arXiv:2406.06367, 2024. 746 Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural Radiance Fields 747 from One or Few Images. In IEEE/CVF Conference on Computer Vision and Pattern Recognition 748 (CVPR), 2021. 749 750 Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun-Juan Zhu, Lionel Ming shuan Ni, and 751 Heung yeung Shum. DINO: DETR with Improved DeNoising Anchor Boxes for End-to-End 752 Object Detection. In The International Conference on Learning Representations (ICLR), 2023. 753 Kai Zhang, Sai Bi, Hao Tan, Yuanbo Xiangli, Nanxuan Zhao, Kalyan Sunkavalli, and Zexiang Xu. 754 GS-LRM: Large Reconstruction Model for 3D Gaussian Splatting. European Conference on 755 Computer Vision, 2024.

756 757 758	Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In <u>IEEE/CVF Conference on Computer</u> <u>Vision and Pattern Recognition (CVPR)</u> , 2018a.
759 760 761 762	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), pp. 586–595, 2018b.
763 764 765	Chuanxia Zheng and Andrea Vedaldi. Free3D: Consistent Novel View Synthesis without 3D Representation. In <u>IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</u> , 2024.
766 767 768 769	Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable DETR: Deformable Transformers for End-to-End Object Detection. In <u>The International Conference on Learning Representations (ICLR)</u> , 2021.
770 771 772 773	Zi-Xin Zou, Zhipeng Yu, Yuan-Chen Guo, Yangguang Li, Ding Liang, Yan-Pei Cao, and Song-Hai Zhang. Triplane Meets Gaussian Splatting: Fast and Generalizable Single-View 3D Reconstruction with Transformers. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 10324–10335, June 2024.
774 775	M. Zwicker, H. Pfister, J. van Baar, and M. Gross. EWA volume splatting. In <u>IEEE Visualization</u> (IEEE VIS), 2001.
776	
777	
778	
779	
780	
781	
782	
783	
784	
785	
786	
787	
788	
789	
790	
(91	
(92 700	
793 704	
94 705	
30 206	
30 207	
J I 00	
90 00	
00	
11	
21 12	
02	
03	
05	
06	
07	
)8	
09	



Figure 6: Spatially Efficient Self-Attention: While employing all queries as query in the selfattention mechanism, we leverage Farthest Point Sampling (FPS) to downsample certain 3D Gaussians. This process enables the extraction of their corresponding queries as keys and values within the self-attention operation.

A APPENDIX

A.1 SPATIAL EFFICIENT SELF ATTENTION (SESA)

While our 3D-aware deformable attention mechanism is notably efficient, the computational cost and memory occupation mainly arises in the self-attention component, particularly when dealing with a large number of 3D Gaussians. However, updating each 3D Gaussian with information from all others is not always necessary because those neighbouring 3D Gaussians usually carry similar information.

To mitigate this issue, as depicted in fig. 6 and drawing inspiration from Wang et al. (2021b), we introduce a technique aimed at reducing the size of the key and value components while maintaining the query component unaltered within the self-attention process. The core concept behind this approach is that while each 3D Gaussian requires updating, not every other 3D Gaussian needs to contribute to this update. We achieve this by selectively updating each query solely with a subset of corresponding queries linked to other 3D Gaussians.

To retain crucial information flow, we leverage the Fast Point Sampling (FPS) algorithm commonly used in point cloud methodologies like PointNet (Qi et al., 2017a) and PointNet++ (Qi et al., 2017b). Specifically, we employ the Gaussian centers μ to identify the most distantly located points and use these points to index the queries. By implementing this strategy, we significantly reduce the model's overall memory footprint while preserving essential information exchange among the Gaussians.

853 854 855

810

811

823

824

825

827

828

829

830

831 832 833

834 835

836

A.2 IMPLEMENTATION DETAILS

856 **Dataset** We utilized a refined subset of the Objaverse LVIS dataset (Deitke et al., 2023) for both training and validating our model. This subset was curated to exclude low-quality models, result-858 ing in a dataset containing 36,044 high-quality objects. This open-category dataset encompasses a 859 diverse range of objects commonly encountered in everyday scenarios. For training, we leveraged 860 rendered images provided by zero-1-to-3 (Liu et al., 2023a) for the random input setting. Each object in the dataset is associated with approximately 12 random views, accompanied by their respective 861 camera poses. We partitioned 99% of the objects for training purposes, reserving the remaining 1% 862 for validation. During training, we randomly selected a subset of views as input while using all 12 863 views for supervision. Each rendered image has a resolution of 512×512 , which we downscaled to



Figure 7: Point clouds of the center of Gaussians from each view. The Gaussians from different views are in different colors.

⁸⁸³ 128×128 . For the fixed view setting, we render the images with fixed views as input and 32 more random views with elevation in (-30, 30) degrees for supervision.

885 To evaluate our model's performance in open-category settings, we conducted tests on the Google 886 Scanned Objects (GSO) benchmark (Downs et al., 2022). The GSO dataset comprises 1,030 3D 887 objects categorized into 17 classes. For this evaluation, we utilized rendered images sourced from Free3D (Zheng & Vedaldi, 2024), which consist of 25 random views along with their corresponding 889 camera poses. Notably, there are no restrictions on the elevation of the rendered views. We utilized 890 the initial views as inputs and the remaining views for assessing our novel view synthesis task. 891 Additionally, we observed that LGM (Tang et al., 2024a) only support fixed-view inputs (e.g., front, 892 left, back, and right). To address this, we evaluated a new rendered GSO dataset at 0 degrees 893 elevation, testing it on 32 random views with elevations ranging from 0 to 30 degrees. To distinguish between the two test sets, we refer to them as GSO-random and GSO-fixed respectively in the 894 following analysis. 895

896

879

880

881 882

Experiment setting We train our model on the setting of 4 views, each time we randomly select 897 4 views as input and all the views for supervision. In the coarse stage, we train the model with less 898 views (i.e. 2 views) with resolution 128×128 and generate 16384 3D Gaussians as initialization 899 of the fine stage. In the fine stage, We use 19600 3D Gaussians to represent the 3D object. For 900 the 3D Gaussians from the coarse stage, we use the mask to remove the background points and 901 padding the number of 3D Gaussians to 19600 by copying some of the remaining 3D Gaussians. 902 The selected 3D Gaussians are then utilized to project queries onto image plane in the refine stage. 903 In each deformable attention layer, we utilize 4 sampling points for each projected 3D Gaussian 904 reference point to sample values on the image.

We use 4 decoder layers and the hidden dimension is 256. Moreover, when training the fine stage, we finetune both the coarse stage and the encoder. We use a mixed-precision training (Narang et al., 2018) with BF16 data type. We train our model with Adam (Kingma & Ba, 2015) optimizer and the learning rate is 0.0001. We take 300K iteration with batch size 4. For the coarse stage, we train it on 8 3090 GPUs (24G) for 5 days and for the fine stage, we train it on 8 A100 (80G) for 3 days.

- 910
- 911 A.3 MORE RESULTS
- 913 A.3.1 QUALITY RESULTS 914

View consistency problem We gives more view inconsistency visualization problem by visualize
 center of Gaussians from each view in different colors, as shown in fig. 7. Gaussians from different
 views representing the same part of the object may lays on the different position in the 3D space and thus cause the view inconsistency problem.





Figure 9: Quality for rendered novel views on GSO-fixed dataset for inputting 4 views with resolution 256 LGM large model.

More visualization We show the point cloud visualization in fig. 8 underscores our model's ability to capture geometry effectively, not just rendering quality.



Figure 10: Quality for rendered novel views on GSO-random dataset for inputting 4 views.

As shown in fig. 9, when given limited number of input, neither LGM nor InstantMesh gives the meanful geomery.

fig. 10 presents the quantitative results of novel views rendered by recent models trained on 4 views.
When provided with 4 random views as input, LGM (Tang et al., 2024a) demonstrates a loss of geometry and encounters view inconsistency problems stemming from its training on fixed views.
In contrast, our approach produces a cohesive 3D Gaussian set that effectively captures object geometries.



Figure 11: Quality for rendered novel views on GSO dataset for inputting 1 view and using Image-Fusion to generate 4 views.

fig. 11 and fig. 12, respectively. The figures illustrate that LGM encounters the issue of view inconsistency; for instance, there are multiple handles visible for the mushroom teapot. InstantMesh loses some details due to its utilization of a discrete triplane to represent continuous 3D space.

fig. 13 shows the result of text-to-3D task. We have incorporated text-to-3D capabilities into our model. To assess quality, we employ MVDream (Shi et al., 2024) to create a single image from a



Figure 12: Quality for rendered novel views on in the wild data for inputting 1 view and using ImageDream to generate 4 views.

Input	NVS					NVS				:	Ga	ussian	center p	oint clo	oud
An astronaut	R		Q.	T	A			Rall	Provide State	R					
A one ear red cup			T			C C		F	1.6	B					
Lion head	6		the second	(E)			C		(F)						
Purple jacket	ß					A	A	M	FIA	<u>p</u>					
Furry dog head	~	C.				4	E		A Star	(S					

Figure 13: Quality for rendered novel views on inputting text and using MVDream to generate 4 views.

text prompt. Subsequently, a diffusion model is utilized to generate multi-view images, which are
 then processed by our model to obtain a 3D representation.

The setting of random input view is obvious a more challenging task than the setting of fixed input view, thus our method also inevitably suffers from a performance drop but still perform better than other state-of-the-art methods. As for Splatter Image (Szymanowicz et al., 2024), it also meets a significant performance drop when random input views are used as its SSIM \uparrow decreased from 0.9151 to 0.8932 and LPIPS \downarrow increased from 0.1517 to 0.2575 despite its PSNR \uparrow has a slight increase. We visualize the results of the two settings to show the difference in fig. 14.

Visualization with resolution 512 We provide the visualization result with resolution 512 in fig. 20.



Single image reconstruction There are common points between our model and TriplaneGaussian and Instant3D that we all use a unitary representation and use Transformer to regress. For



Figure 16: Removing the background use mask for Splatter Image and LGM

1155 Instant3D, it transformers image to Nerf, making longer rendering time. For Triplane Gaussian, 1156 which is a single view reconstruction model with complex and costly triplane representation, rep-1157 resenting compresses 3D space, leading to a lack of detailed information in the 3D structure and 1158 imposing a rigid grid alignment that limits flexibility (Tang et al., 2024a; Qi et al., 2017a). In the contrast, we use a more efficient way (deformable attention) to decode Gaussians. The comparison 1159 between Triplane-Gaussian and our methods is shown in table 5. Triplane Gaussian requires 3D 1160 supervision and takes longer inference time while get worse performance comparing to our model. 1161 We test on the given light-weight checkpoint in the github on the single view situation. We also test 1162 TripoSR (Tochilkin et al., 2024) on the single image reconstruction setting. As shown in table 5, our 1163 model surpass the previous methods on both the performance and the inference speed. We provide 1164 the visualization results of our model on single image reconstruction task in fig. 19 1165

1166 1167

1153 1154

Table 6: Comparison between masked and original pixel aligned methods

Method	$PSNR \uparrow$	SSIM \uparrow	LPIPS \downarrow
LGM	17.4810	0.7829	0.2180
LGM (masked)	21.6008	0.8608	0.1232
Splatter Image	25.6241	0.9151	0.1517
Splatter Image (masked)	25.0648	0.9147	0.1684

1173 1174

1175 **Comparison to masked LGM and Splatter Image** To better explain that the view inconsistency 1176 problem is not caused by the background points from previous methods, we provide the results on 1177 removing background points of LGM and Splatter Image. LGM uses mask loss to make the most 1178 of the pixels contribute to the object itself, even for the background pixels, therefore, removing background use mask makes the results more sparse. It also removing some outliers and thus the 1179 rendering results is better as shown in table 6. Splatter Image keep most of the pixels contribute 1180 to its original position, making most of the background points still located on a plane instead of 1181 the object. Therefore, removing background use mask does not influence the rendering result much 1182 but the rendering quality still reduced a little. Moreover, the view-inconsistency is not caused by 1183 the background points but the mis-alignment of 3D Gaussians from different views, removing the 1184 background use mask does not help solving the problem. We show the visualization in fig. 16 1185

1186

Other number of view results We present the results of training with varying numbers of views (2, 6, 8) and evaluate the corresponding results with the same number of views in table 7.

1188	Table 7: Quantitative results of novel view synthesis training using 2, 6, and 8 input views, tested on
1189	the GSO-random dataset across 2, 6, and 8 views.

Method	$\text{PSNR} \uparrow$	2 views SSIM ↑	LPIPS \downarrow	$PSNR \uparrow$	6 views SSIM ↑	LPIPS \downarrow	$\text{PSNR} \uparrow$	8 views SSIM ↑	LPIPS \downarrow
Splatter Image	22.6390	0.8889	0.1569	26.1225	0.9178	0.1620	26.4588	0.9166	0.1714
Our Model	23.8384	0.8995	0.1254	28.1035	0.9489	0.0559	28.8262	0.9537	0.0492

PSNR of Different Views for Each Method 28 27 26 NSN 522 24 Splatter Image InstantMesh 23 Ours Ż 4 6 Ŕ 10 12 14 16 Number Views

from similar views becomes redundant, so the gain for our model has become plateaued while other methods suffer from performance drop as they cannot handle too many input views due to the view inconsistent problem. As we keep increasing the number of input views larger than 8, our method can still benefit from more input views (as shown in fig. 17) while others meet the CUDA-out-of-memory problem.

Figure 17: Visualization for Splatter Image with fixed view input and random view input.

1209 1210 1211

1197

1198

1199

1201

1202

1203

1205

1207

1208

101

Image-to-3D Image-to-3D conversion represents a fundamental application in 3D genera-

Our model is positioned on the 'sparse view'

setting, which indicates the number of views

less then 10, so we only reports the perfor-

mance of views from 2 to 8 in the main paper. With the increase of input views, information

1212 tion. Following the methodology of LGM and InstantMesh (Tang et al., 2024a; Xu et al., 2024), 1213 we first leverage a multi-view diffusion model, ImageDream (Wang & Shi, 2023), to generate four 1214 predetermined views. Subsequently, our model is employed for 3D Gaussian reconstruction. A 1215 comparative analysis with LGM and InstantMesh is detailed in table 8. For this particular scenario, 1216 we utilize the fixed-view GSO test set with elevations ranging between 0 and 30 degrees. Given 1217 potential variations in camera poses among the generated multi-views, which may not align precisely with standard front, right, back, and left perspectives, we selectively retain 266 objects that 1218 consistently yield accurate images under the provided camera poses. 1219

Table 8: Quantitative results for single view reconstruction on GSO dataset.

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
LGM (Tang et al., 2024a) InstantMesh (Xu et al., 2024)	20.8139 19.4667	0.8581 0.8379	$0.1508 \\ 0.1842$
Our Model	22.3534	0.8567	0.1492

1227 1228

1220

1222

1224 1225 1226

1229 A.4 ABLATION STUDY

Number of views in the coarse stage We add the ablation study on the number of images used during the coarse stage here. The results shown is that the number of images used during the coarse stage does not influence the final result. The reason that we choose the number of views being 2 is that we want to support any number of input views. For example, if we choose the number of views in the coarse stage being 8, we should at least provide 8 views so that the model can not support the number of views smaller than 8. And we tried to change the input views but the number of input views keeping 2 unchanged, the variance of PSNR for 10 different experiments is within 0.185.

1238

1239 Convergence for different regression target Upon investigation, we observe that prior tech-1240 niques frequently predict depth rather than the centers of Gaussians. In our exploration, we con-1241 duct experiments focusing on regressing the centers of 3D Gaussians while keeping other aspects constant. Through this analysis, we discover that regressing the positions of 3D Gaussians can inNumber of views in coarse stage



Table 9: Ablation study results of different view and different number of views for the coarse stage (with 4 views in the refinement stage)

PSNR ↑

SSIM ↑

LPIPS \downarrow

Figure 18: Left: PSNR with different down sampling rate in the spatial efficient self attention.Right: PSNR with different number of Gaussians.

troduce convergence obstacles. Table table 10 illustrates the outcomes of these experiments on the Objaverse validation dataset after 100K steps.

Table 10: Ablation study	on parameter selection.
--------------------------	-------------------------

Coarse-to-fine	25.5338	0.9126	0.0833
Depth 3D Gaussian centers (random initialize in visual cone)	24.3792 19.2551	0.9012 0.8343	0.1014 0.1876
Regression target	$PSNR \uparrow$	SSIM \uparrow	LPIPS \downarrow

1274 1275 1276

1277

1272 1273

1264 1265

1266

1244 1245

More ablation studies Here we gives more ablation study mainly for hyperparameter selection. Due to computational costs, ablation models are trained at 100k iteration and test on Objaverse validation dataset.

1278 1279

1280 **Hyperparameter selection** As previously highlighted, the memory bottleneck of our model lies 1281 in the pointwise self-attention mechanism. To address this, we implement a spatially efficient self-1282 attention technique to alleviate memory consumption. Illustrated in fig. 18 (left), as we augment 1283 the downsampling rate of the key and value in the self-attention mechanism, the memory overhead 1284 diminishes linearly, while the PSNR reduction is not as rapid. Consequently, we opt for a downsampling rate located at the inflection point, which we determine to be 0.01, balancing memory 1285 efficiency with reconstruction quality. Similarly, we select the number of Gaussians as 19600 as 1286 shown in fig. 18 (right). 1287

In table 11, we opted for 4 decoder layers over 6, as the latter offers marginal improvement but demands significantly more computational resources. Additionally, we experimented with using the fine stage initialized with the coarse stage as the encoder and tested the efficacy of fine-tuning both stages. Our findings indicate that fine-tuning both stages yields the best results.

1292

1294

1293 A.5 COMPARISON TO MVSPLAT AND PIXELSPLAT

1295 We present a comparative analysis involving MVSplat (Chen, Yuedong and Xu, Haofei and Zheng, Chuanxia and Zhuang, Bohan and Pollefeys, Marc and Geiger, Andreas and Cham, Tat-Jen and

1296	Table 11: Ablation study on parameter selection.								
1297					F				
1299		Method			PSNR ↑	SSIM \uparrow	LPIPS	↓	
1300		2 decoder layer	rs		24.5229	0.9195	0.1021		
1301		6 decoder layer	rs	adan	26.2442	0.9352	0.0778		
1302		Freeze coarse stage finetune encoder Freeze both coarse stage and encoder			25.3211	0.9223	0.0820		
1303		Default medal	uise stage und ent	couci	26.0212	0.0251	0.0799		
1304					20.2313	0.9551	0.0788		
1305									
1306	Input	1			NVS	1			
1307	mpat	-		_	1115			-	_
1300	4	-	4 4	- 1				4	-
1310	4	1		1				1	1
1311									
1312		_		-	_		_		
1313			- -		3 1				
1314									
1315	×	~			×	-	4	+	1
1316	Kar			1					
1317									
1318	10-5	m		1	🔨 🤺		1	M	The second se
1319	2.2	111	179 1.51	1	()	· /14 9	9	. 9	1.1.
1320									-
1321	<u></u>		e			/		S	
1322	-								
1323	A		A1 A		10 march		and a second		
1324	and the		er 🖉						
1325		:							
1326		Figure 19: Sin	gle view 360 re	ndering	g visualiza	ation on G	SO data	set	
1328									
1329									
1330		Input			N	VS			
1331								-	-
1332									
1333				1			2	T	
1334									
1335					-	-		-	
1336									Ar Car
1337					_				-
1338	- L 🦧 - 🦜	. 15. 🧟	1.00		. 🗶 –	1.10			
1339	₩ ₩₩ ₩₩₩	a (C 🌮 📢 🛓		1		74 🐨 14	6	 // 4%	<i>6</i> 2 t
1340						_			
1342									
1343								2	
1344									
1345	(b) 🐚	()	d 1	50	1	1	- AND		
1346					i in	A Constant	0.5	5 0	CON.
1347						B	H		-
1348									
1349		Figure 20:	Visualization for	or our n	nethod wi	th resolut	ion 512		

1350 Cai, Jianfei, 2024) and pixelSplat (Charatan et al., 2024) on the GSO-random dataset using the 1351 provided checkpoints from the repository in this section. Similar to LGM (Tang et al., 2024a), 1352 both aforementioned methods follow a workflow that regress Gaussians from each views within 1353 the respective camera spaces and subsequently merge them in the global world space. Despite 1354 pixelSplat's integration of cross-view-aware features through an epipolar Transformer, accurately forecasting a dependable probabilistic depth distribution based solely on image features remains 1355 a formidable task (Chen, Yuedong and Xu, Haofei and Zheng, Chuanxia and Zhuang, Bohan and 1356 Pollefeys, Marc and Geiger, Andreas and Cham, Tat-Jen and Cai, Jianfei, 2024). This limitation 1357 often translates to pixelSplat's geometry reconstruction exhibiting comparatively lower quality and 1358 plagued by noticeable noisy artifacts (Chen, Yuedong and Xu, Haofei and Zheng, Chuanxia and 1359 Zhuang, Bohan and Pollefeys, Marc and Geiger, Andreas and Cham, Tat-Jen and Cai, Jianfei, 2024). 1360 Upon examination, we observed that even after isolating points within a visual cone and eliminating 1361 background Gaussians, the geometry fails to convey meaningful information, yielding unsatisfactory 1362 results. 1363

In contrast, MVSplat adopts a design that incorporates a cost volume storing cross-view feature similarities for all possible depth candidates. These similarities offer crucial geometric cues for 3D surface localization, leading to more substantial depth predictions. However, akin to Splatter Image, which assigns each pixel a Gaussian and thereby generates a planar representation rather than the object itself, MVSplat's approach may obscure object details due to occlusion by background Gaussians from other viewpoints, resulting in suboptimal outcomes.

To address this issue, we selectively mask the positioning of Gaussians on background pixels, focusing solely on rendering Gaussians contributing to the object itself. This adjustment reveals significant view inconsistency problems, as illustrated in fig. 21. In the figure, we present the centers of Gaussians generated from different views in different color and the novel views are rendered from the Gaussians from all views. Furthermore, the elaborate incorporation of cross-view attention mechanisms and cost volumes in MVSplat leads to extended inference times and heightened memory requirements as shown in table 12.

Table 12: Comparison with MVSplat and pixelSplat on the GSO-random dataset in the 4-view input setting.

1380	Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Inference time	Rendering time
1381	MVSplat	12.92	0.80	0.30	0.112	0.0090
1382	MVSplat (masked)	16.52	0.80	0.19	0.112	0.0045
1383	pixelSplat (2 views)	12.00	0.80	0.28	1.088	0.0045
1384	pixelSplat (2 views masked)	12.05	0.79	0.27	1.088	0.0023
1385 1386	Ours	26.30	0.93	0.08	0.694	0.0019

1387

1379

1392 1393 1394

1395

1396

1398

1399

1400

1401

1402

