

UniGS: Modeling Unitary 3D Gaussians for Novel View Synthesis from Sparse-view Images

Jiamin Wu^{1,2*}, Kenkun Liu^{3†}, Yukai Shi^{2,4}, Xiaoke Jiang^{2‡}, Yuan YAO¹, Lei Zhang²

¹Hong Kong University of Science and Technology

²International Digital Economy Academy (IDEA)

³The Chinese University of Hong Kong, Shenzhen

⁴Tsinghua University

Abstract

In this work, we introduce **UniGS**, a novel 3D Gaussian reconstruction and novel view synthesis model that predicts a high-fidelity representation of 3D Gaussians from arbitrary number of posed sparse-view images. Previous methods often regress 3D Gaussians locally on a per-pixel basis for each view and then transfer them to world space and merge them through point concatenation. In contrast, Our approach involves modeling unitary 3D Gaussians in world space and updating them layer by layer. To leverage information from multi-view inputs for updating the unitary 3D Gaussians, we develop a DETR (DEtection TRansformer)-like framework, which treats 3D Gaussians as queries and updates their parameters by performing multi-view cross-attention (**MVDEA**) across multiple input images, which are treated as keys and values. This approach effectively avoids ‘ghosting’ issue and allocates more 3D Gaussians to complex regions. Moreover, since the number of 3D Gaussians used as decoder queries is independent of the number of input views, our method allows arbitrary number of multi-view images as input without causing memory explosion or requiring retraining. 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. The code will be released at <https://github.com/jwubz123/UNIG>.

1 Introduction

3D object reconstruction and novel view synthesis (NVS) play a crucial role in the fields of computer vision and graphics. The construction of detailed 3D structures from 2D images has applications including robotics, augmented reality, virtual reality, medical imaging, and beyond. Recently, 3D Gaussian Splatting (3D GS) (Kerbl et al., 2023), as a semi-explicit representation, has demonstrated great efficiency and high-quality rendering performance for NVS.

Despite the advantages of 3D GS, they often perform inadequately when faced with sparse-view inputs (no more than 10 views). Some previous works (Szymanowicz et al., 2024; Tang et al., 2024a; Yi et al., 2024; Zhang et al., 2024b; Xu et al., 2024b) propose to exploit large-scale multi-view image datasets and train feed-forward models that learn to reconstruct 3D Gaussians of objects from sparse-view images in a single forward process. However, many recent feed-forward methods based on 3D GS employ a ‘local prediction and fusion’ approach in which they predict 3D Gaussians for

*This work is done during an internship in the International Digital Economy Academy (IDEA).

†These authors contributed equally to this work.

‡Corresponding author

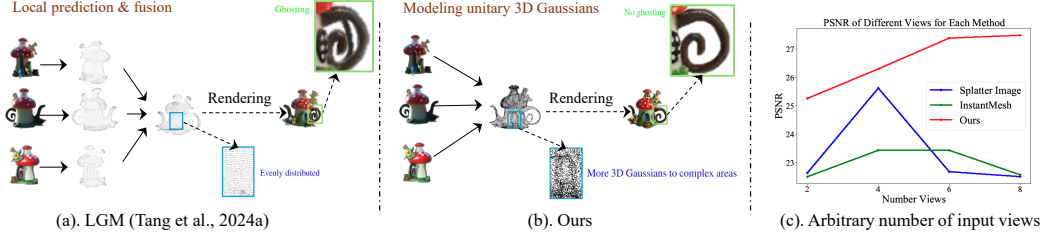


Figure 1: (a) Previous methods like LGM (Tang et al., 2024a) initially predict 3D Gaussians for each pixel for each view and then merge them to get the final 3D Gaussians, resulting in a ‘ghosting’ issue. Moreover, the 3D Gaussians are evenly distributed for both simple and complex regions, while there should be more 3D Gaussians for the complex regions. (b) In contrast, our approach utilizes a unitary set of 3D Gaussians, projecting them onto each view and gathering information across views through a global optimization strategy. Our model effectively avoids the ‘ghosting’ problem and assigns more 3D Gaussians to complex areas (such as the ‘door’ in the image). (c) Our approach supports an arbitrary number of inputs views without requiring retraining and the performance does not deduced. (Trained on 4 input views and tested on 2 to 8 views).

every pixel locally in each input view under the corresponding camera coordinates. Subsequently, the predicted 3D Gaussians are transformed into the same world coordinate system by the input view camera poses and are concatenated to get the final 3D Gaussian representation. This strategy results in even distribution of the positions of 3D Gaussians in both simple and complex regions (Fig. 1) because the local predication on the even distributed pixels on each view. Moreover, the fusion of the 3D Gaussians from each view may encounter the challenge of ‘ghosting’, i.e. the predicted per-pixel 3D Gaussians from different views cannot align to others after concatenation. As depicted in Fig. 1(a), an extra part of the rendered object can be obviously observed. This issue arises from that the estimated positions of 3D Gaussians in each view may be imprecise, so that the merging of 3D Gaussians cannot align well. Another limitation of previous methods is that they require the same number of input views during both training and inference stages; otherwise, their performance would notably deteriorate (Fig. 1(c)), constraining practical applications where the quantity of input views is variable.

To address these issues, we propose a **Unitary 3D Gaussians (UniGS)** representation for sparse-view reconstruction. Inspired by Deformable DETR (DEtection TRansformer) (Zhu et al., 2021) that treats the position and properties of bounding box (Bbox) as queries, we develop a DETR-like framework that treats the unitary 3D Gaussians as queries and updates them layer by layer with multi-view image features as keys and values in cross-attention.

To efficiently leverage multi-view image features, we introduce multi-view deformable cross-attention (**MVDFCA**). Global queries are first defined for the unitary 3D Gaussians within world space. These global queries contribute to generating queries, keys and values for the deformable cross-attention. On the one hand, the global queries undergo modulation by camera parameters for each view to derive view-specific queries, considering that the input contains multiple views but the 3D Gaussians are defined unitarily. Specifically, inspired by (Karras et al., 2019; Hong et al., 2024), the global queries are subject to a linear transformation using weights and biases derived from a multilayer perceptron (MLP) conducted on the camera parameters of each view to obtain the view-specific queries. On the other hand, the global queries are processed through an MLP to predict 3D Gaussians. These 3D Gaussians (centers of them) are projected onto each input view to generate projected points, serving as reference points in deformable DETR within each view. Subsequently, the view-specific queries are updated by deformable cross-attention with keys and values being image features sampled surrounding the reference points for each view. Finally, global queries are updated by all view-specific queries merged through a weighted summation, where the weights are determined by an MLP operating on the view-specific queries.

With the above design, utilizing unitary 3D Gaussians that avoid both per-pixel local optimization and direct fusion of 3D Gaussians from each view, our model effectively addresses the ‘ghosting’ problem, allocates more 3D Gaussians to complex regions, and supports arbitrary number of input views in inference without retraining.

In summary, our contributions are as follows:

- We propose UniGS, a novel 3D object reconstruction and NVS algorithm which introduces a unitary set of 3D Gaussians in world space, enabling all input views to contribute to a unified 3D representation.
- We present MVDFA to efficiently utilize multi-view image features by conducting cross-attention within each view and fuse them into the same set 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 object benchmark GSO.

3D reconstruction from images Recently, various methods have been explored to reconstruct detailed 3D object from limited viewpoints. (Tang et al., 2024a; Liu et al., 2024c; Hong et al., 2024; 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, maintaining view consistency in the generated images is challenging and diffusion models require longer time to generate a single image with denoising process, thus limiting their applicability in real-time scenarios. Moreover, the fusion of information from multiple images remains a non-trivial challenge. Neural Radiance Field (NeRF) (Mildenhall et al., 2020) has gained prominence as a widely used 3D representation (Yu et al., 2021; Cao et al., 2022; Guo et al., 2022; Lin et al., 2022). Techniques such as InstantMesh (Xu et al., 2024a), which combine triplane and NeRF for 3D reconstruction, have demonstrated promising results. (Yu et al., 2021; Wang et al., 2021a; Chen et al., 2021) extract pixel-aligned feature embeddings from multiple views and merge them using MLPs. However, due to its slow rendering speed, NeRF is being supplanted by a new, super-fast, semi-implicit representation—3D Gaussian Splatting (3D GS) (Kerbl et al., 2023). SplatterImage (Szymanowicz et al., 2024), LGM (Tang et al., 2024a) based on 3D Gaussian Splatting, typically handle each input view independently and naively concatenate the resulting 3D Gaussian assets from each view. This method suffers from a lack of information exchange among different views, resulting in inefficient utilization of 3D Gaussians and being view inconsistency. Furthermore, these methods are unable to accommodate an arbitrary number of views as input.

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

2 Related Work

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 process for 3D reconstruction, utilizing Neural Radiance Field (NeRF) (Mildenhall et al., 2020) as a robust 3D representation, have demonstrated effective performance in the field of 3D reconstruction. (Yu et al., 2021; Cao et al., 2022; Guo et al., 2022; Lin et al., 2022; Li et al., 2023b; Müller et al., 2022; Liu et al., 2024d; Wei et al., 2024; Tochilkin et al., 2024; Xu et al., 2024a; Yu et al., 2021; Wang et al., 2021a;

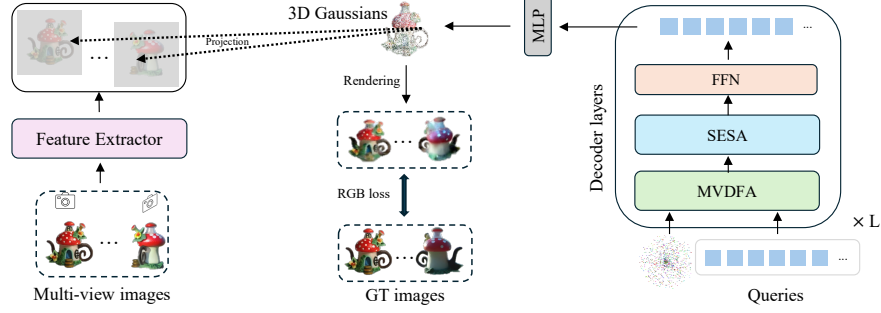
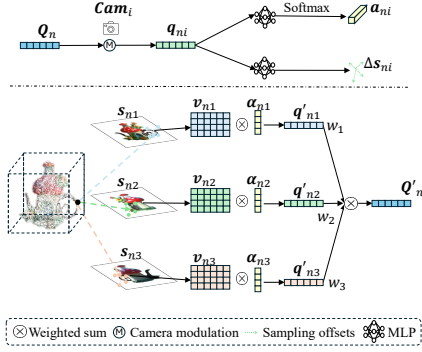


Figure 2: UniGS: Queries are updated by L decoder layers with multi-view image features extracted by the feature extractor. 3D Gaussians are regressed from the queries by an MLP in each layer. Subsequently, they are passed into the next layer and projected onto each view to derive reference points. MVDFA: multi-view deformable attention in Section 3.2. SESA: spatial efficient self-attention in Section 3.4. The dashed arrow means that the centers of 3D Gaussians are projected to multi-view feature maps to retrieve the most related features.

Chen et al., 2021; Liu et al., 2024a; Xiong et al., 2024). 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) (Kerbl et al., 2023). Triplane-Gaussian (Zou et al., 2024), Gamba (Shen et al., 2024), and LeanGaussian (Wu et al., 2025) make promising results on single image 3D reconstruction. When it comes to the parts that does not appear in the input image, the models can not construct them well. Various techniques such as SplatterImage (Szymanowicz et al., 2024), LGM (Tang et al., 2024a), pixelSplat (Charatan et al., 2024), MVSplat (Chen, Yuedong and Xu, Haoifei and Zheng, Chuanxia and Zhuang, Bohan and Pollefeys, Marc and Geiger, Andreas and Cham, Tat-Jen and Cai, Jianfei, 2024), GS-LRM (Zhang et al., 2024b), and GRM (Xu et al., 2024b) have extended the application of 3D Gaussian Splatting to multi-view scenarios. In these approaches, each input view is processed to estimate 3D Gaussians specific to the view, followed by a simple concatenation of the resulting 3D Gaussian assets from all views, resulting in the ‘ghosting’ problem and evenly distribute 3D Gaussians on the object. Such design demands substantial computational resources, particularly as the number of views grows, the number of Gaussians scales linearly with the number of views. Furthermore, these methods are unable to accommodate an arbitrary number of views as input. As an concurrent work, GeoLRM (Zhang et al., 2024a) also exploit the power of transformer and cross-attention, but they still use the discrete voxel representation which may cause high computational cost and potentially get lower accuracy.

Deformable Transformer in 3D Deformable DETR (Zhu et al., 2021) and its following works (Liu et al., 2022; Li et al., 2023a; Liu et al., 2023b; Li et al., 2024; Liu et al., 2023b; Li et al., 2024; Wu et al., 2025; Zhang et al., 2024a, 2023) makes successful attempts in many fields. DFA3D (Li et al., 2023a) and BEVFormer (Li et al., 2022) are introduced to address the feature-lifting challenge in 3D detection and autonomous driving tasks. They achieve notable performance enhancements by employing a deformable Transformer to bridge the gap between 2D and 3D. DFA3D initially uses estimated depth to convert 2D feature maps to 3D, sampling around reference points for deformable attention in each view. However, the 3D sampling point design causes all projected 2D points to represent a singular point, neglecting view variations. BEVFormer (Li et al., 2022) regards the Bird’s-Eye-View (BEV) features as queries, projecting the feature onto each input view. The Spatial Cross-Attention facilitates the fusion of BEV and image spaces, though challenges persist sampling 4 height values per pillar in the BEV feature for selecting 3D reference points may limit coverage, posing challenges in accurate keypoint selection for the model. When contrasting DFA3D and BEVFormer with our MVDFA, a commonality lies in projecting onto 3D regression targets to extract data from various image perspectives. However, our model employ camera modulation to differentiate queries across views, enabling more specific information retrieval.

3 Methods



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

```

1 def MVDDFA(F, K, pi, Q, mu):
2     # Prepare camera embedding
3     camera=concat(K, pi).flatten() #[B, I, 16]
4     cam_embed=MLP(camera) #[B, I, C]
5     # Modulate query by cameras, [B, I, N, C]
6     shift, scale=MLP(cam_embed).chunk(2)
7     q=LayerNorm(Q)*(1+scale)+shift #[B, I, N, C]
8     alpha=softmax(Linear(q)) #Attn score [B, I, N, Ns]
9     # Sampling points
10    Delta s=Linear(q) #[B, I, N, Ns, 2]
11    P=pinhole_proj(camera, mu) #[B, I, N, 1, 2]
12    s=P+Delta s #[B, I, N, Ns, 2]
13    # Weighted sum of view-specific queries
14    V=Linear(F) # [B, I, H*W, C]
15    v=grid_sample(V, s) #[B, I, Ns, C]
16    q'=(alpha*v).sum(-1) #[B, I, N, C]
17    w=sigmoid(Linear(q')) #[B, I, C]
18    Q'=(w*q').sum(-2) #[B, N, C]
19    return Q' #[B, N, C]

```

(b) Pseudo code

Figure 3: MVDDFA: Q_n denotes the n -th unitary queries while q_{ni} denotes the n -th query on the i -th view modulated by the i -th camera Cam_i . Sampling offsets Δs_{ni} and attention score α_{ni} derived by conducting linear transformation on q_{ni} . The sampling offsets are utilized to sample image features at the sampling points $s_{ni} = P_{ni} + \Delta s_{ni}$, where P_{ni} is the reference point derived by projecting the n -th 3D Gaussian. After that, s_{ni} is utilized to sample image features serving as values v_{ni} . These values are employed to update the view-specific queries by attention scores α_{ni} . The unitary queries are refined by the weighted sum of updated view-specific queries q'_{ni} , where w_i is the weight calculated by a linear layer on 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 width W . K and π are camera intrinsics and extrinsics, respectively.

Before introducing details for UniGS, we first give the preliminaries for 3D GS and deformable cross-attention. (see Section 3.1). After that, as illustrated in Fig. 2, our model follows an encoder-decoder framework. All input images undergo processing through an image encoder and a cross-view attention module to extract multi-view image features. For the decoder, 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. 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 (MVDDFA) (see Section 3.2). Spatially efficient self-attention is utilized to reduce computational and memory costs, enabling the utilization of more 3D Gaussians for object reconstruction (see Section 3.4). The training objective is finally introduced. (see Section 3.5).

3.1 Preliminaries

3D GS 3D GS (Kerbl et al., 2023) is a novel rendering method that can be viewed as an extension of point-based rendering methods (Kerbl et al., 2023; Chen & Wang, 2024). They use 3D Gaussians as effective 3D representation for efficient differentiable rendering. The N 3D Gaussians representing an object can be described by $G = \{SH, \mu, \sigma, R, S\}$. The color of 3D Gaussians is represented by spherical harmonics $SH \in \mathbb{R}^{N \times 12}$ while the geometry is described by the center positions $\mu \in \mathbb{R}^{N \times 3}$, shapes (rotation $R \in \mathbb{R}^{N \times 3 \times 3}$ and scales $S \in \mathbb{R}^{N \times 3 \times 4}$), and opacity $\sigma \in \mathbb{R}^{N \times 1}$ of ellipsoids (Zwicker et al., 2001; Kerbl et al., 2023). In our model, the goal for 3D Gaussian reconstruction is to estimate these parameters for a object consists of N 3D Gaussians by a feed-forward network.

Deformable cross-attention Deformable cross-attention, a novel mechanism introduced in Deformable DETR (Zhu et al., 2021), enhances traditional cross-attention by dynamically adjusting spatial sampling locations based on input features. It is initially designed for detection tasks (bounding box center and size prediction). It defines bbox centers as reference points P and its content features

⁴The shape of 3D Gaussian ellipsoids can be described by the covariance matrix Σ , which can be optimized through a combination of rotation and scaling for each ellipsoid as $\Sigma = RSS^T R^T$.

as queries \mathbf{Q} . Keys and values in deformable attention are image features sampled surrounding the reference points at the sampling points \mathbf{s} . To get the sampling points, learnable offsets $\Delta\mathbf{s}$ are predicted by MLP on \mathbf{Q} , and are added to reference points to get the sampling points \mathbf{s} . These learned offsets focus on modeling the most important regions of input feature maps to refine queries. Then the sampling points are utilized to sample image features by grid sampling with bilinear interpolation to get the keys and values of the cross-attention, enabling precise feature attention. The queries are then refined by the image features through deformable attention. Deformable DETR provides a way capturing detailed, context-aware relationships within input data, enhancing overall performance.

3.2 Overview

Feature extractor encoder To extract image features \mathbf{F} from multi-view input, we utilize UNet (Ronneberger et al., 2015; Song et al., 2021b), a widely employed feature extractor in 3D reconstruction tasks, as demonstrated in 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 right after the UNet block, activated when the number of input views exceeds one. In this context, each input view acts as queries, while the concatenation of the remaining views serves as keys and values. To efficiently enable cross-attention across all views, we employ shifted-window attention, as introduced in the Swin Transformer (Liu et al., 2021). This mechanism reduces interactions by focusing on tokens within a local window, effectively reducing memory usage for large input sequences.

Query refinement decoder In the decoder module, we define a fixed number of queries $\mathbf{Q} \in \mathbb{R}^{N \times C}$ with N and C denoting the number of Gaussians and the hidden dimension. The queries are utilized to model 3D Gaussians \mathbf{G} including the center $\boldsymbol{\mu}$, opacity σ , rotation \mathbf{R} , scaling \mathbf{S} , and Spherical Harmonics \mathbf{SH} . Note that queries and 3D Gaussians are one-to-one corresponded. As depicted in Fig. 2, the queries go through multiple decoder layers, each including a multi-view deformable attention (MVDA) (Section 3.3) mechanism to leverage image features, a spatial efficient self-attention (SESA) (Section 3.4) 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, \mathbf{Q}^l is the queries in the l -th layer, and \mathbf{P}^l is reference points in the l -th layer.

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

Finally, queries are processed through an MLP to compute $\Delta\mathbf{G} = \text{MLP}(\mathbf{Q})$ for updating the 3D Gaussian parameters: $\mathbf{G}' = \mathbf{G} + \Delta\mathbf{G}$ ⁵. The initialization for the positions of the center of 3D Gaussians, which make sure that they are not too distant from the ground truth or outside the field of view to guarantee the training convergence, are given in Appendix A.1.2.

3.3 Multi-view deformable attention (MVDA)

View-specific queries generation As depicted Fig. 3 (a), in 3D, the set of queries associated with 3D Gaussian parameters \mathbf{G} are defined unitarily, while cross-attention is defined in multiple image planes. Therefore, the unitary queries should be adjusted to suit each view individually. To solve this issue, we employ 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. More specifically, as shown in Fig. 3(a), the unitary queries \mathbf{Q} are transformed linearly to get the view-specific queries (\mathbf{q}_i for the i -th view) by weights and bias deriving from conducting MLP on camera parameters.

Updating view-specific queries Following deformable DETR (see Section 3.1), the image feature, which are most related to the queries should be sampled surrounding the reference points. These reference points are defined as the position of queries, where in our setting, defined as the projected points in each view for center of 3D Gaussians. More specifically, given center $\boldsymbol{\mu}$ of 3D Gaussians along with the corresponding camera poses $\boldsymbol{\pi}_i$ and intrinsic parameters \mathbf{K}_i for the i -th view, 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 (see Eq. (2)) (Forsyth & Ponce, 2003;

⁵For the updating of 3D Gaussian parameters, rotation is updated by multiplication, while other parameters are updated by addition.

Hartley & Zisserman, 2003),

$$\mathbf{P}_i = \mathbf{K}_i \pi_i \mu \quad (2)$$

. In this context, both matrices \mathbf{K}_i and π_i are expressed in homogeneous form. These projected points \mathbf{P}_i are then regarded as the reference points for 2D deformable attention. After getting the reference points in each view, as shown in Fig. 3, a linear layer is employed to predict the sampling offsets $\Delta \mathbf{s}_i$ to get the sampling points \mathbf{s}_i surrounding the reference points by Eq. (3).

$$\mathbf{s}_i = \mathbf{P}_i + \Delta \mathbf{s}_i \quad (3)$$

Then, we apply the grid sampling algorithm with bilinear interpolation to extract image features at these sampling points, which act as the values \mathbf{v} for cross attention. Subsequently, another linear layer is employed to predicts the attention scores α of the image features \mathbf{v} at the sampling points \mathbf{s} . Finally, for each input view, we compute the updated queries by the dot product between the attention scores α and the sampled values \mathbf{v} .

Fusion of view-specific queries After getting the refined queries for each view \mathbf{q}' , the unitary queries \mathbf{Q}' are then computed as a weighted sum of individual view queries, with the weights calculated using an linear layer on the view-specific queries (see Eq. (4)).

$$\mathbf{Q}' = \sum_i w_i \mathbf{q}'_i, w_i = \text{MLP}(\mathbf{q}'_i) \quad (4)$$

Detailed pseudo code for MVDFA is available in Fig. 3(b).

Insights With MVDFA, 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.

3.4 Spatial efficient self-attention (SESA)

Apart from MVDFA, self-attention in each decoder layer is also important to get the information across all 3D Gaussians. However, when the number of 3D Gaussians N is large, the self-attention is computationally expensive. To tackle this problem, inspired by Wang et al. (2021b), we introduce SESA that reduce the number of keys and values while keeping the number of queries unchanged during self-attention. Insights behind the design is that updating each Gaussian with information from all others may not always be essential, as neighboring Gaussians often contain similar information. This selective updating strategy enables each query to be updated with a subset of related queries instead of all other queries, effectively enhancing the information exchange efficiency. To ensure enough information flow in the downsampled queries, we leverage the Fast Point Sampling (FPS) algorithm from point cloud methodologies (Qi et al., 2017a,b). Specifically, we employ FPS Gaussian centers μ to identify the most distant points from all 3D Gaussians' centers and use the corresponding queries as keys and values in the self-attention. With such strategy, our model optimize memory usage while guaranteeing essential information sharing among Gaussians. Additional details are in Appendix A.1.1.

3.5 Training objective

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 produced by our model. For each object, we render 4 input views and 8 additional views (12 views in total) for supervision. Furthermore, aligning with the methodologies ((Hong et al., 2024; Tang et al., 2024a)), we employ a RGB loss in Eq. (5), which consists of both a mean square error loss \mathcal{L}_{MSE} and a VGG-based LPIPS (Learned Perceptual Image Patch Similarity) loss (Zhang et al., 2018a) $\mathcal{L}_{\text{LPIPS}}$ to guide the rendered views. Here I_{pd} represents the rendered views supervised by the ground truth images I_{gt} . λ is the scalar weiths on the LPIPS loss.

$$\mathcal{L} = \mathcal{L}_{\text{MSE}}(I_{pd}, I_{gt}) + \lambda \mathcal{L}_{\text{LPIPS}}(I_{pd}, I_{gt}) \quad (5)$$

Table 1: Quantitative results for inputting 4 views on GSO dataset. *The results of MV-Gamba and GS-LRM are cited from the paper. ‘NA’ means not reported in the paper. Resolution is 128.

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
Splatter Image (Szymanowicz et al., 2024)	25.6241	0.9151	0.1517
LGM (Small) (Tang et al., 2024a)	17.4810	0.7829	0.2180
LGM (Large) (Tang et al., 2024a)	26.2487	0.9249	0.0541
InstantMesh (Xu et al., 2024a)	23.0177	0.8893	0.0886
GeoLRM* (Zhang et al., 2024a)	22.8400	0.8510	NA
MV-Gamba* (Yi et al., 2024)	26.2500	0.8810	0.0690
GRM (Res-512)* (Xu et al., 2024b)	30.0500	0.9060	0.0520
GS-LRM (Res-256)* (Zhang et al., 2024b)	29.5900	0.9440	0.0510
Our Model	30.4245	0.9614	0.0422

4 Experiments

Dataset We utilize a refined subset of the Objaverse LVIS dataset (Deitke et al., 2023) for training and validation. The training dataset included input rendering images captured from fixed viewpoints (front, back, left, right) and supervised by 32 random views spanning elevations between -30 to 30 degrees. The resolution of the rendered images was downsampled to 128×128 . To evaluate our model, we conducted tests on the Google Scanned Objects (GSO) benchmark 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. More details can be found in Appendix A.2.1.

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 rendered images and the ground truth to evaluate the NVS quality. Additionally, we offer visual representations for both the rendered images and the 3D Gaussian centers as point clouds.

Implementation details are shown in Appendix A.2.2.

4.1 Comparison with state-of-the-art methods

Quantity results We evaluate recent multi-view reconstruction models using 4 views as input. In Table 1, LGM and InstantMesh were evaluated using the provided checkpoints, with “Small” indicating models tailored to 128 resolution with the small model and “Large” to 256 resolution with the large model. Splatter Image (Szymanowicz et al., 2024) is retrained on the same dataset of ours as they do not provide checkpoints with 4 input views. For other methods, we cite the results reported in their paper. Table 1 showcases the performance of these methods in novel view synthesis using 4 fixed views (front, back, right, left) on the GSO 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.3.6.

Quality results We present visualization results for novel view synthesis with resolution 128 in Fig. 4. Note that for LGM, we visualize the results generated by ‘small’ model provided in their github. We can observe the problem of ‘ghosting’ in LGM and the problem of lacking details in InstantMesh (see Appendix A.3.2 for more visualization). Further visualizations with resolution of 256 are accessible in Appendix A.3.3. Furthermore, to demonstrate that our model can handle input views with varying elevations, we also present the results for input views with random camera poses in Appendix A.3.1. The results of comparison to scene-based models (Chen, Yuedong and Xu, Haofer and Zheng, Chuanxia and Zhuang, Bohan and Pollefeys, Marc and Geiger, Andreas and Cham, Tat-Jen and Cai, Jianfei, 2024; Charatan et al., 2024) can be found in Appendix A.3.7. Comparison to the results that removing background points in previous methods are shown in Appendix A.3.5. Results obtained from inputting a single view only can be found in Appendix A.3.4. We also show the point cloud visualization in Fig. 5 underscores our model’s ability to capture geometry effectively, not just rendering quality. More visualization for the ‘ghosting’ problem by visualize center of Gaussians from each view in different colors, are in Appendix A.3.3 Fig. 10. Furthermore, removing the

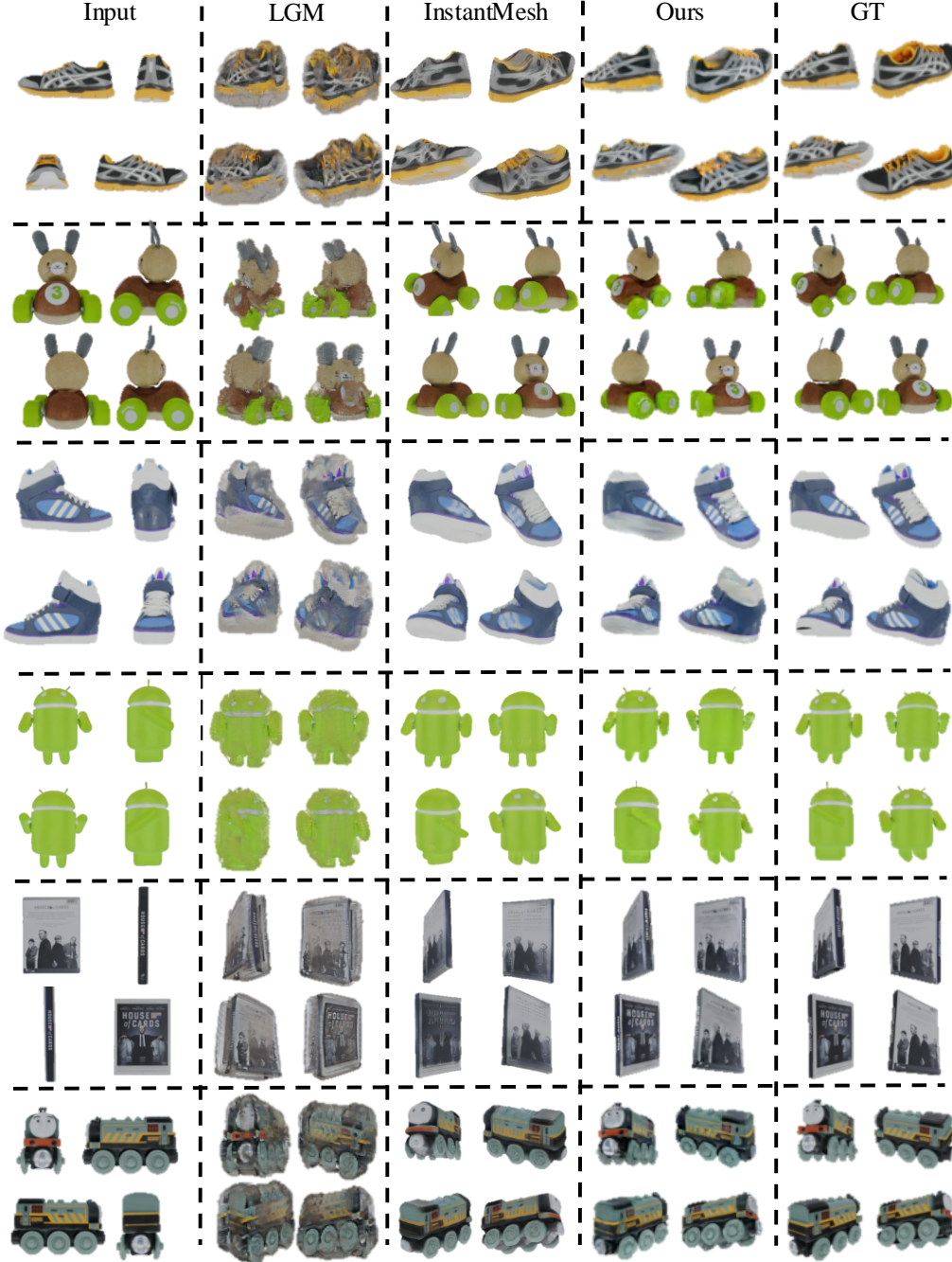


Figure 4: Novel views on GSO dataset for inputting 4 views with resolution 128.

background use masks for Splatter Image and LGM may slightly improve the performance (Fig. 16, Table 7), but still worse than our methods.

Inference on arbitrary number of views Training cost for 3D methods is large, often requiring hundreds of GPUs training multiple days. Additionally, memory costs for previous methods increase linearly with the number of views increasing, presenting challenges for training models with varying input views. Therefore, a model supporting inference with arbitrary number of inputs while being trained on a fixed number of views, such as 4 views, would provide significant advantages. Our model retains unitary 3D Gaussians in world coordinates, treating views as complementary sources

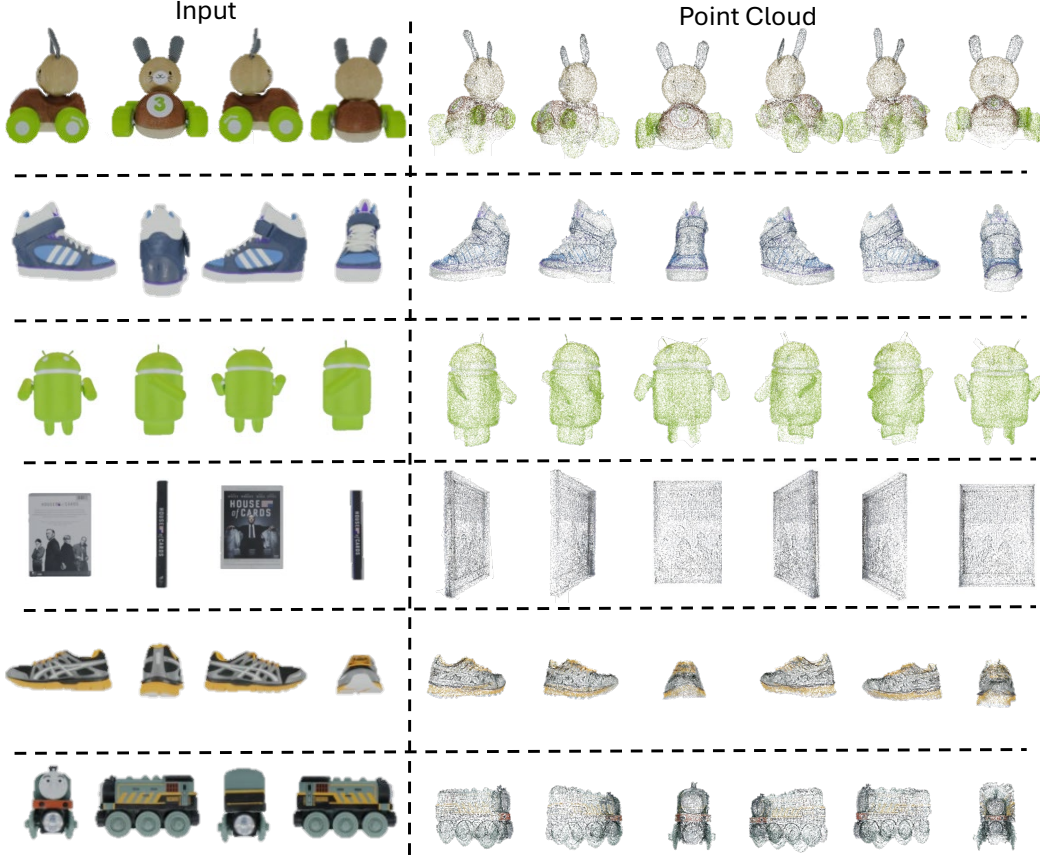


Figure 5: 3D Gaussian center as point cloud on GSO dataset for inputting 4 views.

Table 2: Inference time comparison. 3D: forward time, render: rendering time, inference: time of one forward and 32 rendering. Unit in seconds.

Method	3D ↓	Render ↓	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

without compromising overall 3D integrity. This enables adaptability to variable view counts during inference, despite training on a fixed number of views. Fig. 1 (c) 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.3.6 Fig. 17. While other methods demonstrate satisfactory performance with 4 views during inference, their effectiveness diminishes as the view count different from 4. In contrast, our model gets increasing performance as the number of views increases. It is important that some methods can not handle variations in the number of views between the training and testing, and thus we ignore them in the figure. For the application with single input view, we show the results in Fig. 7 and Table 5 in Appendix A.3.4.

Inference time and memory cost We performed inference time tests across different models, including a diffusion-based method (DreamGaussian (Tang et al., 2024b)), NeRF-based model (InstantMesh (Xu et al., 2024a)), previous Gaussian-based model (LGM (Tang et al., 2024a)), and our model, as shown in Table 2. In contrast to previous methods that compute 3D Gaussians per pixel per input view, our model retains a single 3D Gaussian irrespective of the number of views.

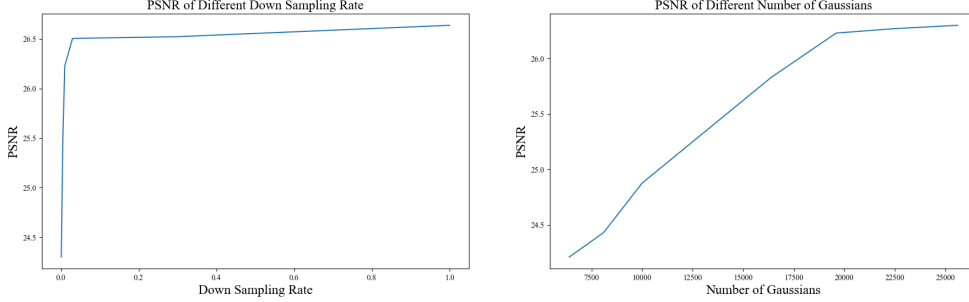


Figure 6: Left: PSNR with different down sampling rate in SESA. Right: PSNR with different number of Gaussians.

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.

4.2 Ablation studies

Table 3 illustrates an ablation study that evaluates different components of the model architecture. All the experiments are evaluated on the Objaverse validation dataset.

Initialization of 3D Gaussians As shown in Table 3, initialize 3D Gaussians randomly (without any constraint) results in low performance. This problem arises from utilizing image features around the projected 3D Gaussian center within each image view. When projections extend beyond the image plane, the reference points may be out of the image plane and sampled images features with only value 0, which steep gradients. Initialize 3D Gaussians randomly within the cone of vision (CoV) make the performance better. In this sense, better initialization may give better final performance and therefore we initialize the 3D Gaussians by regressing them first from image features for each pixel. Details for such initialization are provided in Appendix A.1.2.

No cross-view attention in feature extractor Removing cross-view attention leads to a moderate decrease in performance compared to the full model.

No decoder We do ablation study of removing the decoder, i.e., use the 3D Gaussians regressed from each pixel of the image features extracted from the feature extractor and then concatenate them as the same process of previous methods (Tang et al., 2024a; Szymanowicz et al., 2024; Zhang et al., 2024b). As shown in Appendix A.1.2, such design underperforms the full model.

No camera modulation in MVDA Furthermore, removing the camera modulation on queries or use the shared 3D sampling points instead of sampling on each view adversely impacts the results, underscoring the critical significance of this view-specific design. The full model achieves the best performance, indicating that each component contributes positively to the overall model effectiveness.

Ablation study on SESA As introduced in Section 3.4, the memory bottleneck of our model lies in the pointwise self-attention mechanism. To address this, we implement a spatially efficient self-attention technique to alleviate memory consumption. Illustrated in Fig. 6 (left), as we augment the downsampling rate of the key and value in the self-attention mechanism, the memory overhead diminishes linearly, while the PSNR reduction is not so rapid. Consequently, we opt for a downsampling rate located at the inflection point, which we determine to be 0.01, balancing memory efficiency with reconstruction quality.

Number of Gaussians In our model, we first define a fixed number of 3D Gaussians. To show influence of the selection on the number of 3D Gaussians, we conduct ablation study on different

Table 3: Ablation study on model design.

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
ran. init.	12.1213	0.6531	0.6224
ran. init. in CoV	22.6740	0.8711	0.2383
w/o cross view attention	25.3923	0.9013	0.1007
UNet only	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

number of 3D Gaussians. As shown in Fig. 6 (right), when the number of 3D Gaussians exceeds 19600, the increment from adding the number of Gaussians becomes flatten. Therefore, we select the number of 3D Gaussians being 19600. Additionally, we offer details on hyperparameter selections in Appendix A.4.

Number of input views We add the ablation study on the number of input views in Appendix A.4 Table 10. Our model not only support 4 views as input, but also surpass previous methods on other number of input views. More ablation studies are in Appendix A.4.

4.3 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., 2024a), 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.3.3. We also showcase the quality results of our model on the GSO dataset in Fig. 7. Additionally, we also provide the results of in-the-wild input images in Appendix A.3.3 Fig. 11.

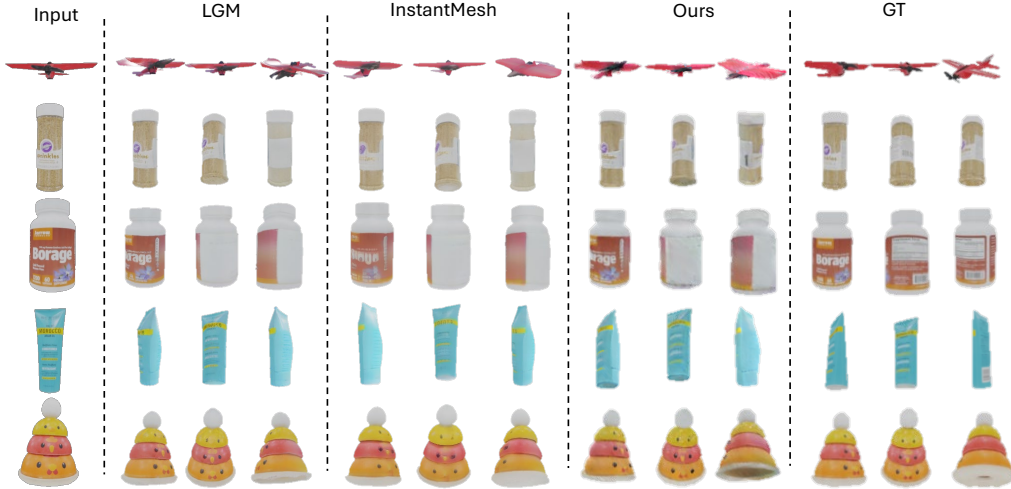


Figure 7: Quality for rendered novel views on GSO dataset for inputting 1 view and using ImageFusion to generate 4 views.

Text-to-3D 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.3.3.

5 Conclusion and Limitation

In this paper, we introduce a novel sparse view 3D Gaussian reconstruction and NVS method. First, fixed number of unitary 3D Gaussians defined in world space together with the corresponding queries are initialized, and each 3D Gaussian (center of them) is projected onto input image features extracted by a feature extractor. Then, MVDFA block is designed to do cross-attention on each view. The unitary queries are modulated to each view by the camera parameters to get the view-specific queries. The 3D Gaussians (center of them) are projected onto multi-view image plane to get the reference points. After that, image features surrounding these reference points are sampled to update the view-specific queries by deformable cross-attention in each view. Finally, the unitary queries are refined by weighted sum on the view-specific queries. Moreover, we develop a spatially efficient self-attention mechanism to minimize computational costs. With the above design, our model successfully tackling ‘ghosting’ problem and giving more meaningful distribution of 3D Gaussians by giving more 3D Gaussians for complex areas. Moreover, our model can accept an arbitrary number of views as input without damage the performance. The experiments on GSO dataset shows the effectiveness of our model.

Limitations Our model requires user-provided camera parameters, which are important for projection, presenting potential challenges in 3D reconstruction.

References

- 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.
- 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. ArXiv, 2024.
- Chen, Yuedong and Xu, Haofei and Zheng, Chuanxia and Zhuang, Bohan and Pollefeys, Marc and Geiger, Andreas and Cham, Tat-Jen and Cai, Jianfei. Mvsplat: Efficient 3d gaussian splatting from sparse multi-view images. European Conference on Computer Vision, 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 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 13142–13153, 2023.
- 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 2022 International Conference on Robotics and Automation (ICRA), pp. 2553–2560. IEEE, 2022.
- David A Forsyth and Jean Ponce. A Modern Approach. Computer vision: a modern approach, 17: 21–48, 2003.
- 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. Multiple View Geometry in Computer Vision. 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.
- Tero Karras, Samuli Laine, and Timo Aila. A Style-Based Generator Architecture for Generative Adversarial Networks. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR), pp. 4401–4410, 2019.
- Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing 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.
- Bernhard Kerbl, Georgios Kopanas, Thomas Leimkuehler, and George Drettakis. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. In ACM Transactions on Graphics (TOG), 2023.
- Diederik Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In International Conference on Learning Representations (ICLR), San Diego, CA, USA, 2015.
- Hongyang Li, Hao Zhang, Zhaoyang Zeng, Shilong Liu, Feng Li, Tianhe Ren, and Lei Zhang. DFA3D: 3D Deformable Attention For 2D-to-3D Feature Lifting. In Proceedings of the IEEE/CVF international conference on computer vision (ICCV), 2023a.

- Hongyang Li, Hao Zhang, Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, and Lei Zhang. TAPTR: Tracking Any Point with Transformers as Detection. In Proceedings of the IEEE/CVF European Conference on Computer Vision (ECCV), 2024.
- Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan Sunkavalli, Greg Shakhnarovich, and Sai Bi. Instant3D: Fast Text-to-3D with Sparse-View Generation and Large Reconstruction Model. The International Conference on Learning Representations (ICLR), 2023b.
- Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai. BEVFormer: Learning Bird’s-Eye-View Representation from Multi-Camera Images via Spatiotemporal Transformers. European conference on computer vision (ECCV), 2022.
- Kai-En Lin, Yen-Chen Lin, Wei-Sheng Lai, Tsung-Yi Lin, Yichang Shih, and Ravi Ramamoorthi. Vision Transformer for NeRF-Based View Synthesis from a Single Input Image. In 2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), 2022.
- Kenkun Liu, Derong Jin, Ailing Zeng, Xiaoguang Han, and Lei Zhang. A Comprehensive Benchmark for Neural Human Radiance Fields. Advances in Neural Information Processing Systems (NIPS), 36, 2024a.
- Minghua Liu, Ruoxi Shi, Linghao Chen, Zhuoyang Zhang, Chao Xu, Xinyue Wei, Hansheng Chen, Chong Zeng, Jiayuan Gu, and Hao Su. One-2-3-45++: Fast Single Image to 3D Objects with Consistent Multi-View Generation and 3D Diffusion. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10072–10083, 2024b.
- Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Mukund Varma T, Zexiang Xu, and Hao Su. One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. In Conference on Neural Information Processing Systems (NIPS), 2024c.
- Minghua Liu, Chong Zeng, Xinyue Wei, Ruoxi Shi, Linghao Chen, Chao Xu, Mengqi Zhang, Zhaoning Wang, Xiaoshuai Zhang, Isabella Liu, Hongzhi Wu, and Hao Su. MeshFormer: High-Quality Mesh Generation with 3D-Guided Reconstruction Model. Conference on Neural Information Processing Systems (NIPS), 2024d.
- Ruoshi Liu, Rundui Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot One Image to 3D Object. In IEEE/CVF International Conference on Computer Vision (ICCV), 2023a.
- Shilong Liu, Feng Li, Hao Zhang, Xiao Yang, Xianbiao Qi, Hang Su, Jun Zhu, and Lei Zhang. DAB-DETR: Dynamic anchor boxes are better queries for DETR. In International Conference on Learning Representations (ICLR), 2022.
- 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 IEEE/CVF International Conference on Computer Vision (ICCV), 2023b.
- 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 Proceedings of the IEEE/CVF international conference on computer vision (ICCV), pp. 10012–10022, 2021.
- 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 The European Conference on Computer Vision (ECCV), 2020.
- Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. In ACM Transactions on Graphics (SIGGRAPH), 2022.
- 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 6th international conference on learning representations (ICLR), volume 1, pp. 14, 2018.

- Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation. In Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR), pp. 652–660, 2017a.
- Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space. Advances in neural information processing systems (NIPS), 30, 2017b.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR), pp. 10684–10695, 2022.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation. In Medical image computing and computer-assisted intervention (MICCAI), 2015.
- Qihong 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. arXiv preprint arXiv:2403.18795, 2024.
- Yichun Shi, Peng Wang, Jianglong Ye, Long Mai, Kejie Li, and Xiao Yang. MVDream: Multi-view Diffusion for 3D Generation. The International Conference on Learning Representations (ICLR), 2024.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising Diffusion Implicit Models. In International Conference on Learning Representations (ICLR), 2021a.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-Based Generative Modeling through Stochastic Differential Equations. The International Conference on Learning Representations (ICLR), 2021b. URL <https://openreview.net/forum?id=PxTIG12RRHS>.
- Stanislaw Szymanowicz, Christian Rupprecht, and Andrea Vedaldi. Splatter Image: Ultra-Fast Single-View 3D Reconstruction. In IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2024.
- 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.
- Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. DreamGaussian: Generative Gaussian Splatting for Efficient 3D Content Creation. In International Conference on Learning Representations (ICLR), 2024b.
- Dmitry Tochilkin, David Pankratz, Zexiang Liu, Zixuan Huang, , Adam Letts, Yangguang Li, Ding Liang, Christian Laforte, Varun Jampani, and Yan-Pei Cao. TripoSR: Fast 3D Object Reconstruction from a Single Image. arXiv preprint arXiv:2403.02151, 2024.
- Yuri Viazovetskyi, Vladimir Ivashkin, and Evgeny Kashin. StyleGAN2 Distillation for Feed-forward Image Manipulation. In Proceedings of the IEEE/CVF European Conference on Computer Vision (ECCV), pp. 170–186. Springer, 2020.
- Peng Wang and Yichun Shi. ImageDream: Image-Prompt Multi-view Diffusion for 3D Generation. arXiv preprint arXiv:2312.02201, 2023.
- Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. IBRNet: Learning Multi-View Image-Based Rendering. In 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Jun 2021a. doi: 10.1109/cvpr46437.2021.00466. URL <http://dx.doi.org/10.1109/cvpr46437.2021.00466>.
- Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In Proceedings of the IEEE/CVF international conference on computer vision (ICCV), pp. 568–578, 2021b.

- Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image Quality Assessment: From Error Visibility to Structural Similarity. IEEE transactions on image processing, 13(4): 600–612, 2004.
- 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 preprint arXiv:2404.12385, 2024.
- Jiamin Wu, Kenkun Liu, Han Gao, Xiaoke Jiang, and Lei Zhang. LeanGaussian: Breaking Pixel or Point Cloud Correspondence in Modeling 3D Gaussians. Conference on Computer Vision and Pattern Recognition (CVPR), 2025.
- Zhangyang Xiong, Chenghong Li, Kenkun Liu, Hongjie Liao, Jianqiao Hu, Junyi Zhu, Shuliang 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 Computer Vision and Pattern Recognition (CVPR), pp. 19801–19811, 2024.
- Jiale Xu, Weihao Cheng, Yiming Gao, Xintao Wang, Shenghua Gao, and Ying Shan. InstantMesh: Efficient 3D Mesh Generation from a Single Image with Sparse-view Large Reconstruction Models. arXiv preprint arXiv:2404.07191, 2024a.
- Yinghao Xu, Zifan Shi, Wang Yifan, Sida Peng, Ceyuan Yang, Yujun Shen, and Wetzstein Gordon. GRM: Large Gaussian Reconstruction Model for Efficient 3D Reconstruction and Generation, author=Xu Yinghao and Shi Zifan and Yifan Wang and Chen Hansheng and Yang Ceyuan and Peng Sida and Shen Yujun and Wetzstein Gordon, 2024b.
- Xuanyu Yi, Zike Wu, QiuHong Shen, Qingshan Xu, Pan Zhou, Joo-Hwee Lim, Shuicheng Yan, Xinchao Wang, and Hanwang Zhang. MVGamba: Unify 3D Content Generation as State Space Sequence Modeling. arXiv preprint arXiv:2406.06367, 2024.
- Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural Radiance Fields from One or Few Images. In IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021.
- Chubin Zhang, Hongliang Song, Yi Wei, Yu Chen, Jiwen Lu, and Yansong Tang. GeoLRM: Geometry-Aware Large Reconstruction Model for High-Quality 3D Gaussian Generation. NeurIPS, 2024a.
- Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun-Juan Zhu, Lionel Ming shuan Ni, and Heung yeung Shum. DINO: DETR with Improved DeNoising Anchor Boxes for End-to-End Object Detection. In The International Conference on Learning Representations (ICLR), 2023.
- Kai Zhang, Sai Bi, Hao Tan, Yuanbo Xiangli, Nanxuan Zhao, Kalyan Sunkavalli, and Zexiang Xu. GS-LRM: Large Reconstruction Model for 3D Gaussian Splatting. European Conference on Computer Vision, 2024b.
- Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2018a.
- 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.
- Chuanxia Zheng and Andrea Vedaldi. Free3D: Consistent Novel View Synthesis without 3D Representation. In IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2024.
- Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable DETR: Deformable Transformers for End-to-End Object Detection. In The International Conference on Learning Representations (ICLR), 2021.
- 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.

M. Zwicker, H. Pfister, J. van Baar, and M. Gross. EWA volume splatting. In IEEE Visualization (IEEE VIS), 2001.

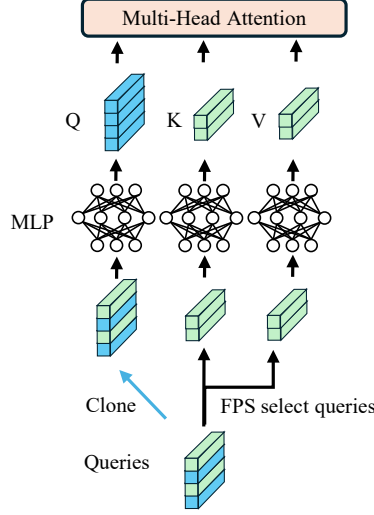


Figure 8: Spatially Efficient Self-Attention: While employing all queries as query in the self-attention mechanism, we leverage Farthest Point Sampling (FPS) to downsample certain 3D Gaussians. This process enables the extraction of their corresponding queries as keys and values within the self-attention operation.

A Appendix

A.1 Model details

A.1.1 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. 8 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 insights 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.

More specifically, 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). 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, our model significantly reduces the model’s overall memory footprint while preserving essential information exchange among the Gaussians.

A.1.2 Analysis for the regression of the center of 3D Gaussians and the initialization for 3D Gaussians and queries

From camera space to world space In Szymanowicz et al. (2024), the Gaussian centers are located in each input view’s camera space as shown in Eq. (6),

$$\mu_{cam} = \begin{bmatrix} x_{cam} \\ y_{cam} \\ z_{cam} \end{bmatrix} = \begin{bmatrix} u_1 d + \Delta_x \\ u_2 d + \Delta_y \\ d + \Delta_z \end{bmatrix} \quad (6)$$

The center coordinates $x_{cam}, y_{cam}, z_{cam}$ are parameterized by depth d and offsets $(\Delta_x, \Delta_y, \Delta_z)$. u_1, u_2 are the UV coordinates of the ray passing through the corresponding input image. This design

represents each point with multiple Gaussians, potentially introducing ‘ghosting’ due to concatenation issues at various points caused by depth inaccuracies and tend to shortcut input views (Wu et al., 2025). In our framework, we define unitary Gaussians in world space, project their centers to each input view for feature retrieval, as depicted in Fig. 3(a). The centers of Gaussians can be written as $\mu_{\text{world}} = [x_{\text{world}}, y_{\text{world}}, z_{\text{world}}]$.

3D Gaussian initialization However, during the initial training phases, discrepancies between the 3D Gaussian centers and ground truth often result in imprecise selection of image features at sampling points, presenting challenges for model convergence as shown in Section 4.2. To address this issue, we utilize a initialization that directly regress 3D Gaussian parameters from each pixel of the image features. We first train a coarse network that predicts 3D Gaussians for each pixel of each input view under the camera space of that view. After that, the 3D Gaussians are transformed and fused to get the 3D Gaussians under world space. The role of this network is to provide a coarse initialization of 3D Gaussians for the subsequent refinement. We use the UNet architecture, which has the same structure to the feature extractor.

Moreover, 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 the reference view. This approach allows us to represent all 3D data within this consistent relative coordinate system.

A.2 Experiment details

A.2.1 Datasets

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, resulting in a dataset containing 36,044 high-quality objects. This open-category dataset encompasses a diverse range of objects commonly encountered in everyday scenarios. For training, we leveraged 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 camera poses. We partitioned 99% of the objects for training purposes, reserving the remaining 1% for validation. During training, we randomly selected a subset of views as input while using all 12 views for supervision. Each rendered image has a resolution of 512×512 , which we downsampled to 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.

To evaluate our model’s performance in open-category settings, we conducted tests on the Google Scanned Objects (GSO) benchmark (Downs et al., 2022). The GSO dataset comprises 1,030 3D 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 camera poses. Notably, there are no restrictions on the elevation of the rendered views. We utilized the initial views as inputs and the remaining views for assessing our novel view synthesis task. Additionally, we observed that LGM (Tang et al., 2024a) only support fixed-view inputs (e.g., front, left, back, and right). To address this, we evaluated a new rendered GSO dataset at 0 degrees 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 and GSO respectively in the following analysis.

A.2.2 Implementation details

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

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

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
Splatter Image	25.7660	0.8932	0.2575
LGM	15.1113	0.8440	0.1592
InstantMesh	17.3073	0.8525	0.1376
Our Model	26.3020	0.9255	0.0836

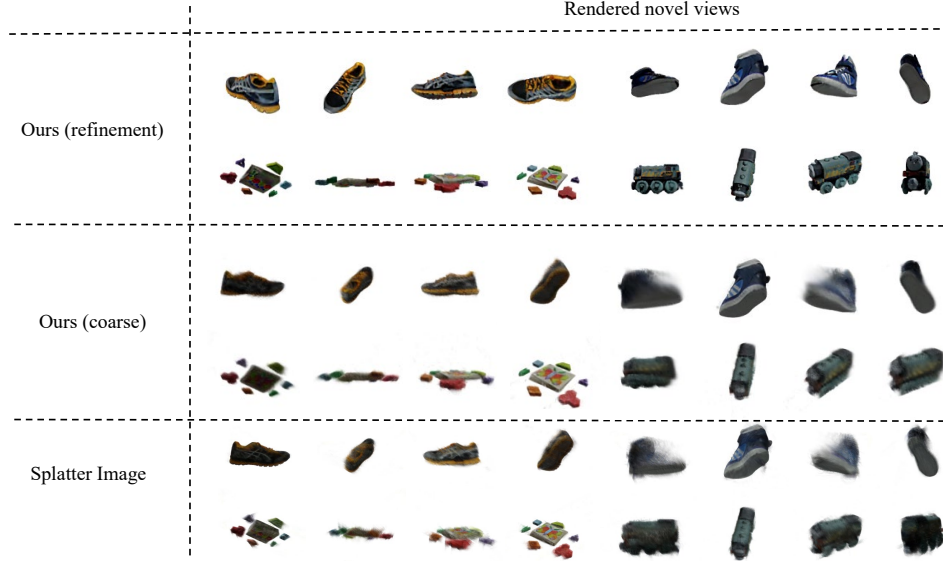


Figure 9: Visualization for Splatter Image with fixed view input and random view input.

We use 4 decoder layers and the hidden dimension is 256. 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 initialization, we train it on 8 3090 GPUs (24G) for 5 days and for the refinement stage, we train it on 8 A100 (80G) for 3 days.

A.3 More results

A.3.1 Input views with random camera poses

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 4 displays the results when utilizing random 4 views as input on the GSO dataset. Notably, there is a performance drop observed in LGM and InstantMesh with random input views. For Splatter Image, although the PSNR does not reduced much, its SSIM and LPIPS reduced significantly. We provide more visualization in Fig. 15.

PSNR of Splatter Image in Table 4 is good but SSIM and LPIPS are not good enough, we further provide the visualization in Fig. 9.

A.3.2 More example for ‘Ghosting’ problem

We gives more ‘ghosting’ visualization problem by visualize center of Gaussians from each view in different colors, as shown in Fig. 10. 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 ‘ghosting’ problem.

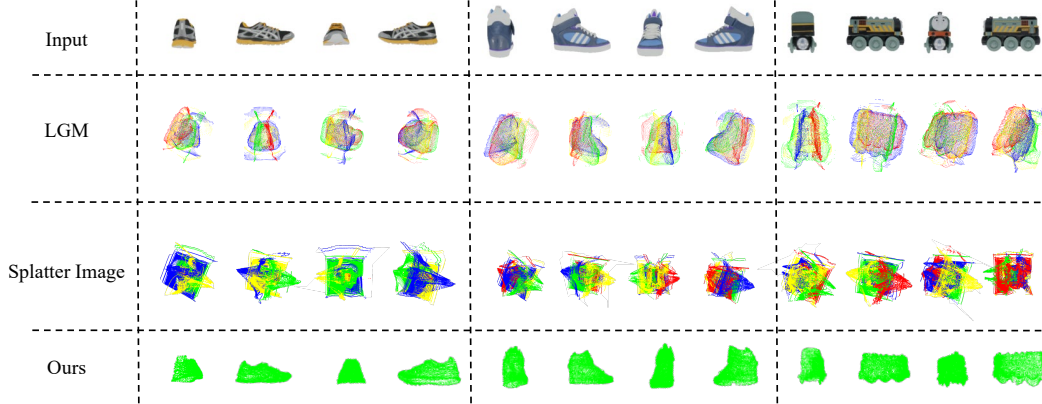


Figure 10: Point clouds of the center of Gaussians from each view. The Gaussians from different views are in different colors.

A.3.3 More visualization

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., 2024a), we first leverage a multi-view diffusion model, ImageDream (Wang & Shi, 2023), to generate four predetermined views. Subsequently, our model is employed for 3D Gaussian reconstruction. A comparative analysis with LGM and InstantMesh is detailed in Table 5. For this particular scenario, we utilize the fixed-view GSO test set with elevations ranging between 0 and 30 degrees. Given 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 consistently yield accurate images under the provided camera poses.

Table 5: Quantitative results for single view reconstruction on GSO dataset.

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

As shown in Fig. 12, when given limited number of input, neither LGM nor InstantMesh gives the meaningful geomery.

Fig. 13 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 ‘ghosting’ problems stemming from its training on fixed views. In contrast, our approach produces a cohesive 3D Gaussian set that effectively captures object geometries.

The figures illustrate that LGM encounters the issue of ‘ghosting’; 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. 14 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 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

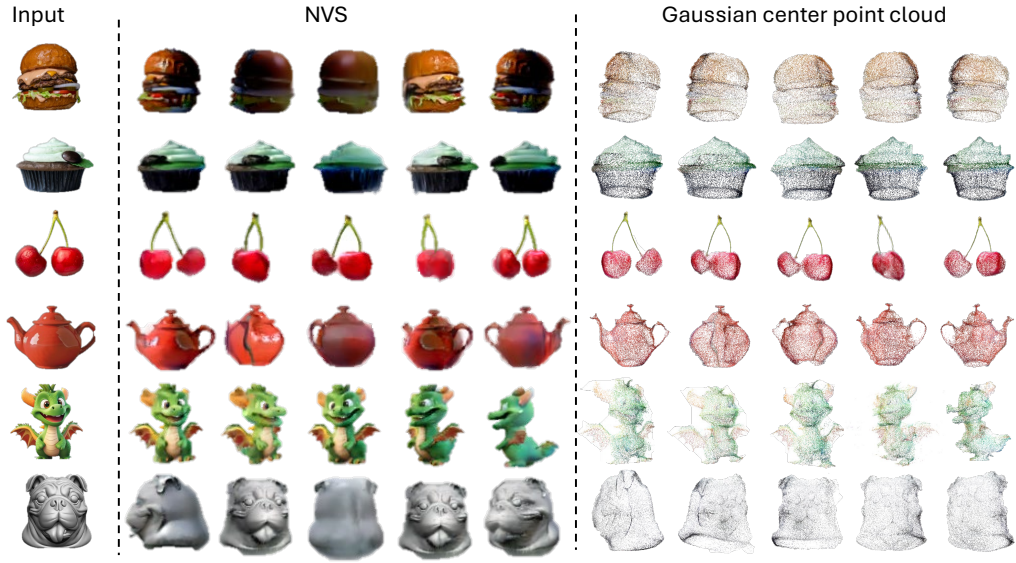


Figure 11: Quality for rendered novel views on in the wild data for inputting 1 view and using ImageDream to generate 4 views.

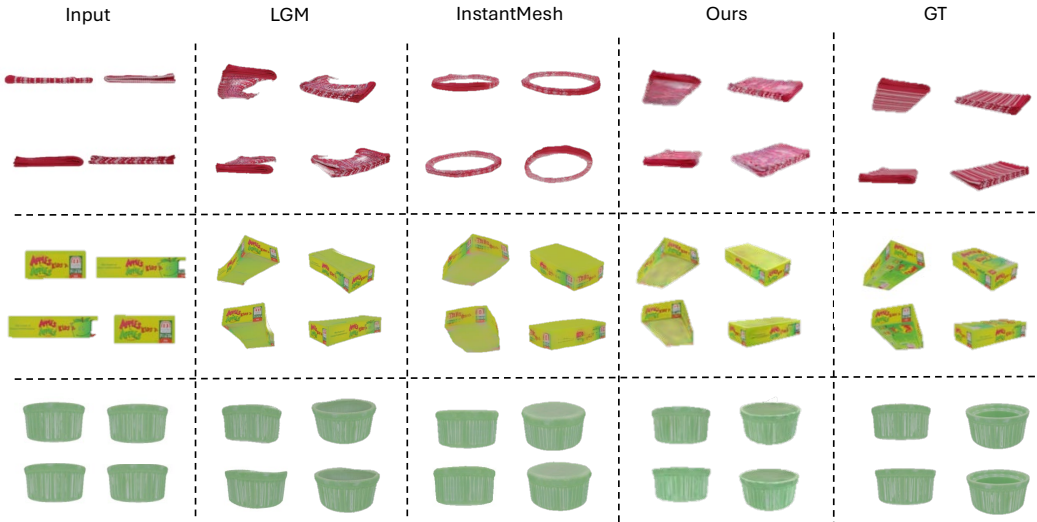


Figure 12: Quality for rendered novel views on GSO dataset for inputting 4 views with resolution 256 LGM large model.

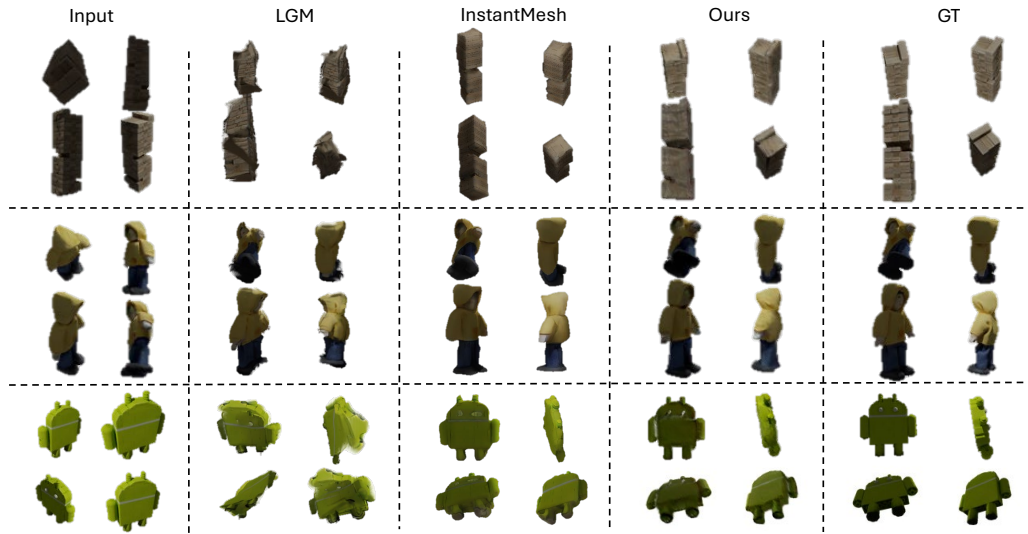


Figure 13: Quality for rendered novel views on GSO dataset for inputting 4 views with random camera poses.

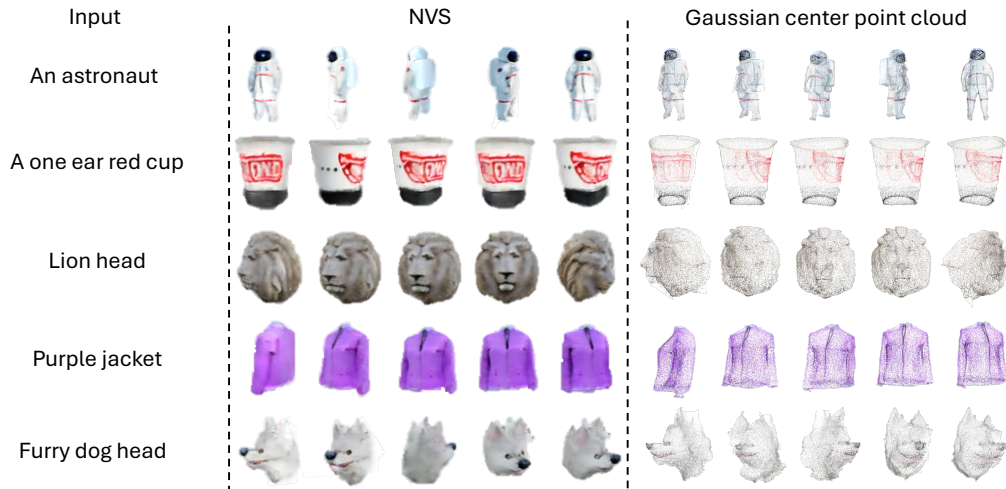


Figure 14: Quality for rendered novel views on inputting text and using MVDream to generate 4 views.

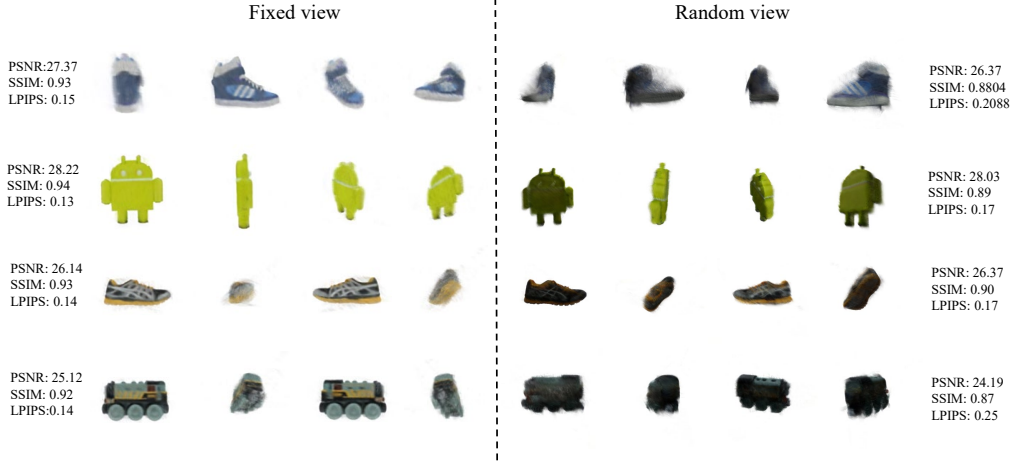


Figure 15: Visualization for Splatter Image with fixed view input and random view input.

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. 15.

We provide the visualization result with resolution 512 in Fig. 20.

A.3.4 Single image reconstruction

Table 6: Quantitative results trained on Objaverse LVIS and tested on GSO. 3D sup. means need 3D supervision.

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	3D sup.	Inference time	Rendering time
Triplane-Gaussian (Zou et al., 2024)	18.61	0.853	0.159	✓	1.906	0.0025
TripoSR (Tochilkin et al., 2024)	20.00	0.872	0.149	✗	3.291	22.7312
Ours	23.45	0.897	0.093	✗	0.476	0.0025

There are common points between our model and TriplaneGaussian and Instant3D that we all use a unitary representation and use Transformer to regress. For Instant3D, it transformers image to Nerf, making longer rendering time. For Triplane Gaussian, which is a single view reconstruction model with complex and costly triplane representation, representing compresses 3D space, leading to a lack of detailed information in the 3D structure and 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 between Triplane-Gaussian and our methods is shown in Table 6. Triplane Gaussian requires 3D supervision and takes longer inference time while get worse performance comparing to our model. We test on the given light-weight checkpoint in the github on the single view situation. We also test TripoSR (Tochilkin et al., 2024) on the single image reconstruction setting. As shown in Table 6, our model surpass the previous methods on both the performance and the inference speed. We provide the visualization results of our model on single image reconstruction task in Fig. 19

A.3.5 Comparison to masked LGM and Splatter Image

Table 7: 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

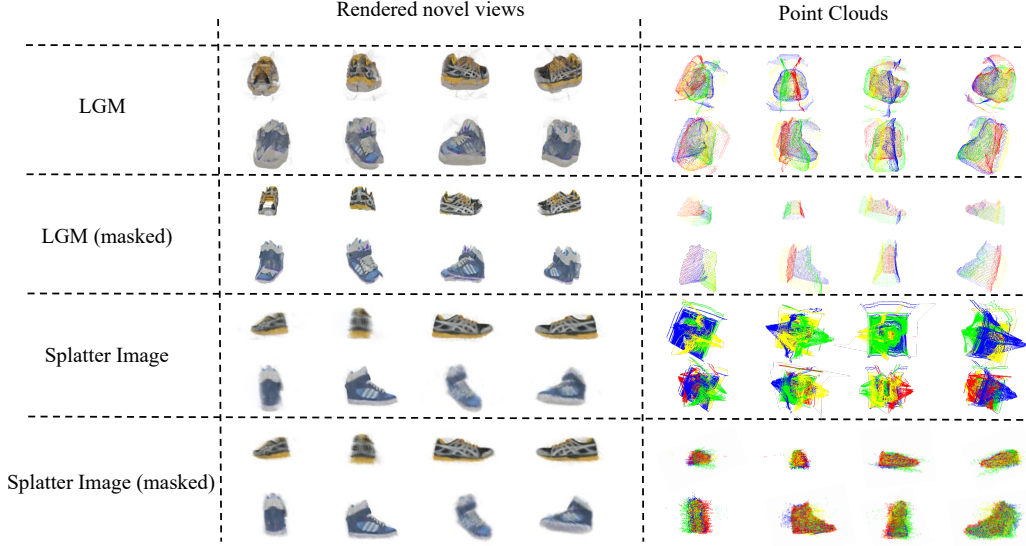


Figure 16: Removing the background use mask for Splatter Image and LGM

To better explain that the ‘ghosting’ problem is not caused by the background points from previous methods, we provide the results on removing background points of LGM and Splatter Image. LGM uses mask loss to make the most 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 rendering results is better as shown in Table 7. Splatter Image keep most of the pixels contribute to its original position, making most of the background points still located on a plane instead of the object. Therefore, removing background use mask does not influence the rendering result much but the rendering quality still reduced a little. Moreover, the ‘ghosting’ is not caused by the background points but the mis-alignment of 3D Gaussians from different views, removing the background use mask does not help solving the problem. We show the visualization in Fig. 16

A.3.6 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 8.

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

Method	2 views			6 views			8 views		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	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

Our model is positioned on the ‘sparse view’ setting, which indicates the number of views less than 10, so we only reports the performance of views from 2 to 8 in the main paper. With the increase of input views, information 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.

A.3.7 Comparison to MVSpLat and pixelSpLat

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 Cai,

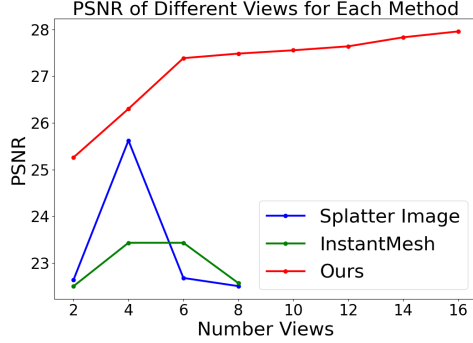


Figure 17: Visualization for Splatter Image with fixed view input and random view input.

Jianfei, 2024) and pixelSplat (Charatan et al., 2024) on the GSO dataset by training it on Objaverse. Similar to LGM (Tang et al., 2024a), both aforementioned methods follow a workflow that regress Gaussians from each views within the respective camera spaces and subsequently merge them in the global world space. Despite 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 a formidable task (Chen, Yuedong and Xu, Haoifei and Zheng, Chuanxia and Zhuang, Bohan and Pollefeys, Marc and Geiger, Andreas and Cham, Tat-Jen and Cai, Jianfei, 2024). This limitation often translates to pixelSplat’s geometry reconstruction exhibiting comparatively lower quality and plagued by noticeable noisy artifacts (Chen, Yuedong and Xu, Haoifei and Zheng, Chuanxia and Zhuang, Bohan and Pollefeys, Marc and Geiger, Andreas and Cham, Tat-Jen and Cai, Jianfei, 2024). Upon examination, we observed that even after isolating points within a visual cone and eliminating background Gaussians, the geometry fails to convey meaningful information, yielding unsatisfactory results.

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 ‘ghosting’ problems, as illustrated in Fig. 18. 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 9.

Table 9: Comparison with MVSplat on the GSO dataset in the 4-view input setting.

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Inference time	Rendering time
MVSplat	23.06	0.90	0.13	0.112	0.0090
MVSplat (masked)	24.10	0.91	0.12	0.112	0.0045
Ours	26.30	0.93	0.08	0.694	0.0019

A.4 Ablation study

Number of views for the initialization We add the ablation study on the number of images used during the the training of initialization. The results shown is that the number of images used 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 being 8, we



Figure 18: Visualization for MVSplat and our method

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.

Table 10: Ablation study results of different view and different number of views for the initialization (with 4 views in the refinement stage)

Number of views	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
1	30.2312	0.9608	0.0413
2	30.4245	0.9614	0.0422
3	30.3442	0.9618	0.0419
4	30.4521	0.9620	0.0412

Convergence for different regression target Upon investigation, we observe that prior techniques frequently predict depth rather than the centers of Gaussians. In our exploration, we conduct 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 introduce convergence obstacles. Table Table 11 illustrates the outcomes of these experiments on the Objaverse validation dataset after 100K steps.

Table 11: Ablation study on parameter selection.

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

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.

Table 12: Ablation study on parameter selection.

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
2 decoder layers	24.5229	0.9195	0.1021
6 decoder layers	26.2442	0.9352	0.0778
Freeze encoder	25.3211	0.9264	0.1003
Default model	26.2313	0.9351	0.0788

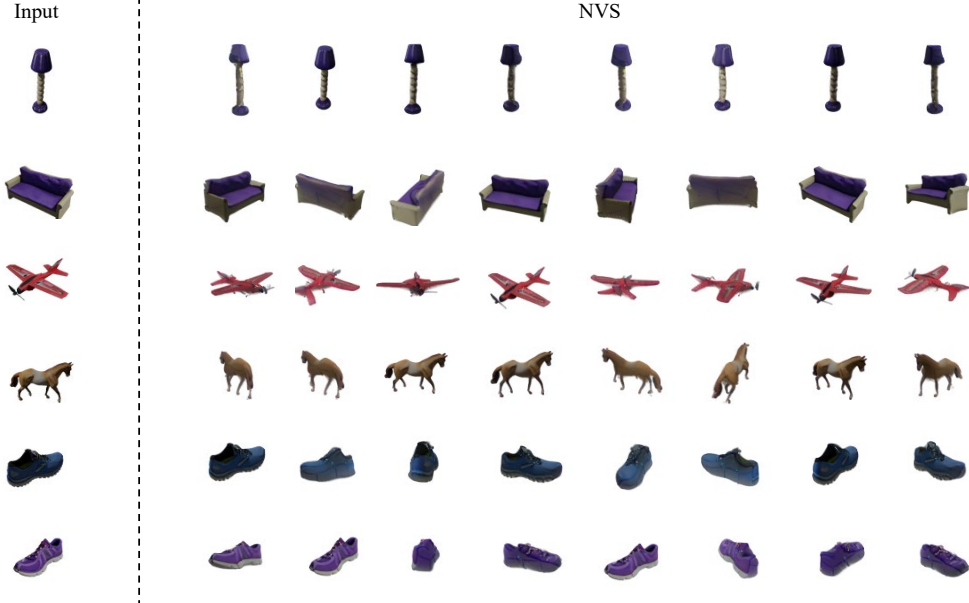


Figure 19: Single view 360 rendering visualization on GSO dataset

Hyperparameter selection In Table 12, we opted for 4 decoder layers over 6, as the latter offers marginal improvement but demands significantly more computational resources. Our findings indicate that fine-tuning yields the better results.

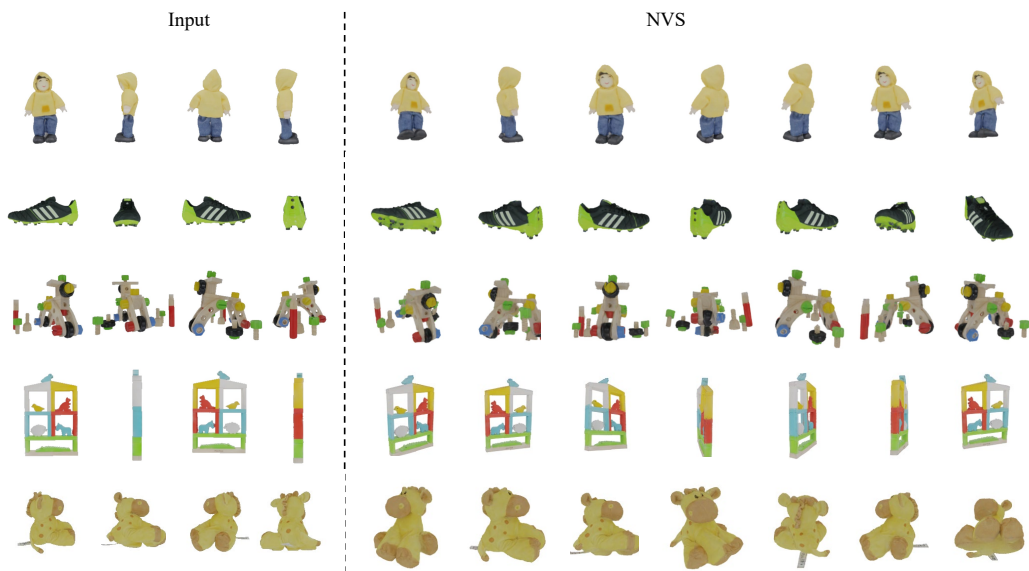


Figure 20: Visualization for our method with resolution 512