# Gamba: Marry Gaussian Splatting with Mamba for Single-View 3D Reconstruction

**Qiuhong Shen**[1*]  **Zike Wu**[3*]  **Xuanyu Yi**[3*]  **Pan Zhou**[2,4†]  **Hanwang Zhang**[3,5]

**Shuicheng Yan**[5]  **Xinchao Wang**[1†]

[1]National University of Singapore  [2]Singapore Management University
[3]Nanyang Technological University  [4]Sea AI Lab  [5]Skywork AI

## Abstract

We tackle the challenge of efficiently reconstructing a 3D asset from a single image at millisecond speed. Existing methods for single-image 3D reconstruction are primarily based on Score Distillation Sampling (SDS) with Neural 3D representations. Despite promising results, these approaches encounter practical limitations due to lengthy optimizations and significant memory consumption. In this work, we introduce Gamba, an end-to-end 3D reconstruction model from a single-view image, emphasizing two main insights: (1) Efficient Backbone Design: introducing a Mamba-based GambaFormer network to model 3D Gaussian Splatting (3DGS) reconstruction as sequential prediction with linear scalability of token length, thereby accommodating a substantial number of Gaussians; (2) Robust Gaussian Constraints: deriving radial mask constraints from multi-view masks to eliminate the need for warmup supervision of 3D point clouds in training. We trained Gamba on Objaverse and assessed it against existing optimization-based and feed-forward 3D reconstruction approaches on the GSO Dataset, among which Gamba is the only end-to-end trained single-view reconstruction model with 3DGS. Experimental results demonstrate its competitive generation capabilities both qualitatively and quantitatively and highlight its remarkable speed: Gamba completes reconstruction within 0.05 seconds on a single NVIDIA A100 GPU, which is about $1,000\times$ faster than optimization-based methods. Please see our project page at https://florinshen.github.io/gamba-project.

## 1 Introduction

We tackle the challenge of efficiently reconstructing a 3D asset from a single image, an endeavor with substantial implications across diverse industrial sectors. This endeavor facilitates AR/VR content generation from a single snapshot and aids in the development of autonomous vehicle path planning through monocular perception [44, 15, 62].

Previous approaches to single-view 3D reconstruction have mainly been achieved through Score Distillation Sampling (SDS) [37], which leverages pre-trained 2D diffusion models [8, 40] to guide optimization of the underlying representations of 3D assets. These optimization-based approaches have achieved remarkable success, known for their high-fidelity and generalizability. However, they require a time-consuming per-instance optimization process [46, 56, 59] to generate a single object and also suffer from artifacts such as the "multi-face" problem arising from bias in pre-trained 2D diffusion models [17]. On the other hand, previous approaches predominantly utilized neural
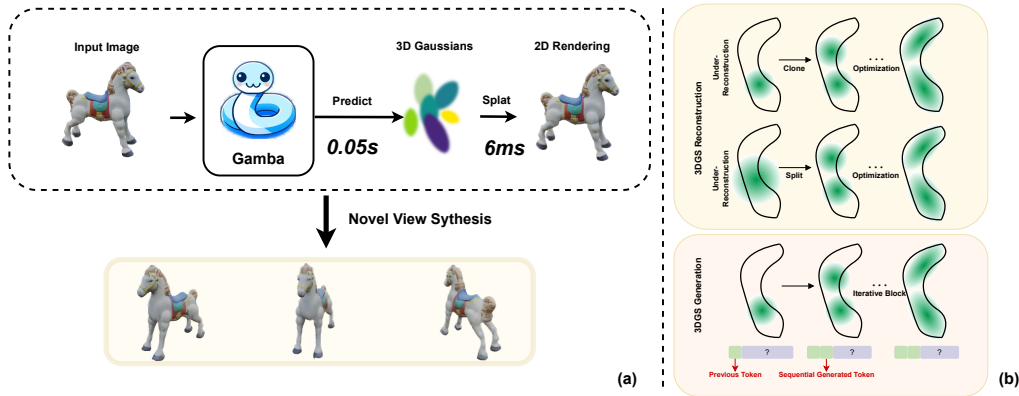
---

Figure 1: (a): We propose Gamba, an end-to-end, feed-forward single-view reconstruction pipeline, which marries 3D Gaussian Splatting with Mamba to achieve fast reconstruction. (b): The relationship between the 3DGS iterative reconstruction and the Gamba sequential prediction pattern.

radiance fields (NeRF) [32, 1], which are equipped with high-dimensional multi-layer perception (MLP) and inefficient volume rendering [32]. This computational complexity significantly limits practical applications on limited compute budgets. For instance, the Large reconstruction Model (LRM) [18] is confined to a resolution of 32 using a triplane-NeRF [43] representation, and the resolution of renderings is limited to 128 due to the bottleneck of online volume rendering.

To address these challenges and thus achieve *efficient* single-view 3D reconstruction, we are seeking an amortized reconstruction framework with the groundbreaking 3D Gaussian Splatting, notable for its memory-efficient and high-fidelity tiled rendering [20, 69, 3, 55]. Despite recent exciting progress [48], how to properly and immediately generate 3D Gaussians remains a less studied topic. Recent prevalent 3D amortized generative models [18, 53, 60, 61, 68, 22] predominantly use transformer-based architecture as their backbones [49, 36], but we argue that these widely used architectures are sub-optimal for generating 3DGS. The crucial challenge stems from the fact that 3DGS requires a sufficient number of 3D Gaussians to accurately represent a single 3D object. However, the computational complexity of the attention mechanism in transformers increases quadratically with the number of tokens [49]. Existing works [60, 68] struggle to build amortized 3DGS reconstruction models with transformers. They resort to a multi-stage training paradigm to alleviate this bottleneck, but the performance is still limited due to the insufficient number of Gaussians in the first stage. Furthermore, 3DGS has explicit, non-structural, and discrete nature, making the simultaneous generation of 3DGS parameters a more challenging task compared to its Neural Radiance Fields (NeRF) counterparts.

To tackle the above challenges, we start by revisiting the iterative 3DGS reconstruction from posed multi-view images. The analysis presented in Figure 1(b) reveals that the densification during the iterative 3DGS reconstruction can be conceptualized as a sequential prediction based on previously predicted tokens. With this insight, we introduce a novel architecture for *end-to-end* 3DGS reconstruction dubbed Gaussian Mamba (Gamba), which is built upon a new scalable sequential network, Mamba [9]. Our Gamba enables context-dependent reasoning and scales linearly with sequence (token) length, allowing it to efficiently mimic the inherent process of 3DGS reconstruction when reconstructing 3D assets enriched with a sufficient number of 3D Gaussians. Moreover, we train Gamba on the large-scale Objaverse dataset [5] with a robust training constraints in an end-to-end manner. At its core, the model employs a radial mask constraint to supervise the placement of Gaussians effectively, thereby eliminating the need for an explicit point cloud supervision and multi-stage training in previous work [68, 60]. Due to its feed-forward, end-to-end architecture, combined with efficient rendering of 3DGS, Gamba achieves remarkable speed. It requires only about 0.05 seconds to generate a 3D asset and 6 ms for synthesizing novel views, which is $1000\times$ faster than previous optimization-based methods [57, 38] still delivering comparable quality in reconstruction outputs.

We demonstrate the superiority of Gamba on the wide range of single images and decent evaluation on the Google Scanned Object (GSO) dataset [6]. Both qualitative and quantitative experiments clearly indicate that Gamba can instantly generate high-quality and diverse 3D assets from a single image, continuously outperforming other state-of-the-art methods. In summary, we make three-fold contributions:

- We introduce GambaFormer, a simple Mamba-based reconstructor to process 3D Gaussian Splatting, which has global context length with linear complexity.

- Integrated with GambaFormer and robust 3DGS constraint, we present Gamba, an end-to-end 3DGS reconstruction pipeline for efficient single-view reconstruction.

- Extensive experiments show that Gamba outperforms the state-of-the-art baselines in terms of reconstruction quality and speed.

## 2   Related Works

**Amortized 3D Generation.** Amortized 3D generation is able to instantly generate 3D assets in a feed-forward manner after training on large-scale 3D datasets [58, 5, 64], in contrast to tedious SDS-based optimization methods [59, 25, 57, 16, 46]. Previous works [34, 33] married de-noising diffusion models with various 3D explicit representations (*e.g.*, point cloud and mesh), which suffers from lack of generalizablity and low texture quality. Recently, pioneered by LRM [18], several works utilize the capacity and scalability of the transformer [36] and propose a full transformer-based regression model to decode a NeRF representation from triplane features. The following works extend LRM to predict multi-view images [22], combine with diffusion [61], and pose estimation [53]. However, their triplane-NeRF representation is restricted to inefficient volume rendering and relatively low resolution with blurred textures. Gamba instead seeks to train an efficient feed-forward model marrying Gaussian splatting with Mamba for single-view 3D reconstruction.

**Gaussian Splatting for 3D Generation.** The explicit nature of 3DGS facilitates real-time rendering capabilities and unprecedented levels of control and editability, making it highly relevant for 3D generation. Several works have effectively utilized 3DGS in conjunction with optimization-based 3D generation [59, 37, 25]. For example, DreamGaussian [48] utilizes 3D Gaussian as an efficient 3D representation that supports real-time high-resolution rendering via rasterization. Despite the acceleration achieved, generating high-fidelity 3D Gaussians using such optimization-based methods still requires several minutes and a large computational memory demand. TriplaneGaussian [68] extends the LRM architecture with a hybrid triplane-Gaussian representation. AGG [60] decomposes the geometry and texture generation task to produce coarse 3D Gaussians, further improving its fidelity through Gaussian Super Resolution. Splatter image [45] and PixelSplat [2] propose to predict 3D Gaussians as pixels on the output feature map of two-view images. LGM [47] generates high-resolution 3D Gaussians by fusing information from multi-view images generated by existing multi-view diffusion models [42, 52] with an asymmetric U-Net. Among them, our Gamba demonstrates its superiority and structural elegance with *single image* as input and an *end-to-end*, *single-stage*, feed-forward manner.

## 3   Method

In this section, we detail our proposed single-view 3D reconstruction pipeline, which incorporates 3D Gaussian Splatting (3DGS) as depicted in Figure 2(a), dubbed as "Gamba." The core component of this pipeline is the GambaFormer, which predicts 3D Gaussians from a single image input (see Sec. 3.2). We design elaborate constraints on the Gaussian parameters and a progressive training strategy, as discussed in Sec. 3.3, to achieve end-to-end training and high-fidelity reconstruction.

### 3.1   Preliminary of 3D Gaussian Splatting

3D Gaussian Splatting (3DGS) [20] has gained prominence as an efficient explicit 3D representation, using anisotropic 3D Gaussians to achieve intricate modeling. Each Gaussian is defined by its 3D central position $\mu \in \mathbb{R}^3$, covariance matrix $\Sigma$, associated color $c \in \mathbb{R}^3$ (applicable when the degree of spherical harmonics is set to zero), and opacity $\alpha \in \mathbb{R}$. To be better optimized, the covariance matrix $\Sigma$ is constructed from a 3D scale $r \in \mathbb{R}^3$ and a rotation quaternion $q \in \mathbb{R}^4$. Generally, the $j$-th Gaussian can be collectively denoted as $\mathcal{G}_j = \{\mu_j, \alpha_j, r_j, q_j, c_j\}$. 3DGS projects Gaussians onto 2D images using a tile-based rasterization pipeline to support real-time rendering and differentiable optimization. This approach effectively controls the number of Gaussians through both adaptive densification and pruning of Gaussians.

camera pose token
condition image tokens
3DGS tokens

Image Tokenizer

Gamba Block $L_0$
Gamba Block $L_1$
Gamba Block $L_{N-1}$
3DGS Decoder

3D GS Render

3DGS params

Cam poses

Novel views supervision

$\mathbf{G}_n$
Drop
Mamba Block
Prepend
Linear
$\mathbf{G}_{n-1}$

a) Overall architecture of Gamba
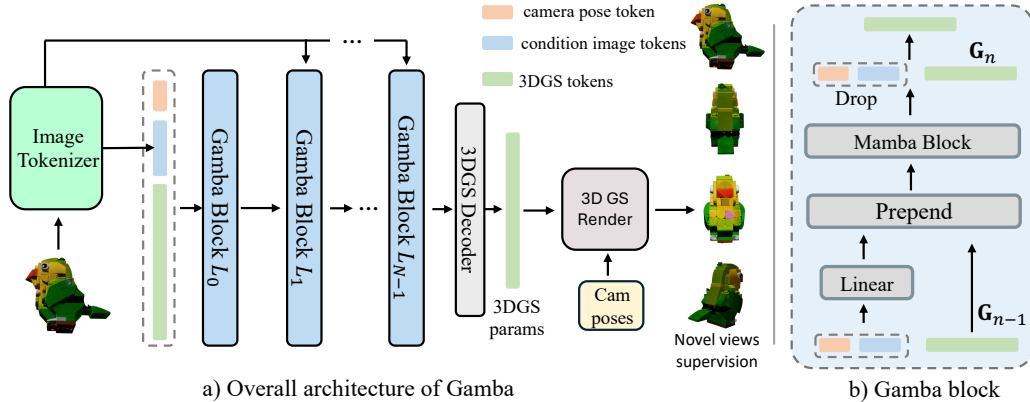
b) Gamba block

Figure 2: Overall architecture of Gamba. Gamba takes a single view image and its camera pose as input to predict the 3D Gaussian Splatting of the given subject. Training supervision is only applied on the rendered multi-view images through reconstruction loss.

## 3.2 GambaFormer

Given a set of multi-view images and their paired camera pose $\{\mathbf{x}_i, \pi_i\}$ of a 3D object, Gamba first transforms the reference image $\mathbf{x}_{\text{ref}}$ and pose $\pi_{\text{ref}}$ into condition tokens. These tokens are then concatenated with the learnable 3DGS tokens to predict a set of 3D Gaussians. Subsequently, the predicted Gaussians are rendered into 2D multi-view images using the given camera poses $\{\pi\}$. This rendering process employs the differentiable rasterizer of 3DGS [20], enabling direct supervision of the multi-view rendering output by the provided ground-truth images $\{\mathbf{x}\}$ at both reference and novel viewpoints through image space reconstruction loss.

**Condition image tokens.** The reference view $\mathbf{x} \in \mathbb{R}^{H \times W \times 3}$ is tokenized using the visual transformer (ViT) model DINO [35], which has demonstrated robust feature extraction capabilities in previous large reconstruction models (LRM) [18, 68]. This process extracts the reference image $\mathbf{x}$ into a sequence of tokens $\mathbf{X} \in \mathbb{R}^{K \times C}$, characterized by a length of $K$ and channel dimensions $C$.

**Condition camera tokens.** Given the variability of camera poses $\pi_{\text{ref}}$ across the sampled 3D objects during the training phase, it is essential to embed camera features as a condition in our GambaFormer. Following the precedent settings in LRM [18], we construct the camera matrix using 12 parameters that include the rotation and translation components of the camera extrinsics and 4 parameters $[fx, fy, cx, cy]$ representing the camera intrinsics. These parameters are then transformed into a high-dimensional camera embedding $\mathbf{T} \in \mathbb{R}^C$ via a multi-layer perceptron (MLP). It is important to note that Gamba does not rely on any canonical pose, and the ground truth $\pi$ is required solely as input during training for multi-view supervision.

**Expanding image as 3DGS tokens.** To effectively reconstruct a 3D object using a set of Gaussians, our framework necessitates the prediction of a substantial number of 3D Gaussians. Achieving a sufficient count of Gaussians is essential for accurately fitting a 3D object. Previous methods [68, 60] have resorted to a two-stage training framework to manage the considerable memory overhead associated with long token sequences, which initially trained a network with capability of predicting up to $L = 4096$ Gaussians, and then trained a super-resolution network enhances the resolution of the output from the first stage to $L = 16384$.

In contrast, our GambaFormer architecture obviates this two-stage training paradigm by leveraging the linear complexity of state space models. The Gaussian count is set as $L = 16384$ throughout our framework. To construct the 3DGS token sequence, we start by embedding the paired camera pose $\pi_{ref}$ into the reference image $\mathbf{x}_{ref}$ with Plücker rays [61] at each image pixel, which are then concatenated into $\mathbf{s} \in \mathbb{R}^{H \times W \times 9}$. Following this, a large non-overlapping convolution with a kernel size of $p \times p$ is applied to transform $\mathbf{s}$ into a feature map $\mathbf{S} \in \mathbb{R}^{h \times w \times D}$:

$$\mathbf{G} = \text{Scan}(\text{Conv}(\mathbf{s})) + \mathbf{E}, \tag{1}$$

where $h = H/p$ and $w = W/p$. The dimension of each Gamba block is denoted by $D$. We then employ four pre-defined scan orders [28] to flatten this feature map into 1D sequence of length

4

$L = 4 \times h \times w$. Finally, this sequence is plus with learnable 3D Gaussian Splatting (3DGS) embeddings, $\mathbf{E} \in \mathbb{R}^{L \times D}$, resulting in the formation of 3DGS tokens, $\mathbf{G} \in \mathbb{R}^{L \times D}$, which serve key input to our GambaFormer.

**Core of the Gamba Block.** The detailed architecture of the Gamba block, compared with the vanilla Mamba block, is illustrated in Figure 2(b). While the Mamba block excels at processing long sequences of tokens, existing variants [28, 67, 23] have not explored traditional cross-attention mechanisms. Leveraging the unidirectional scan order inherent to Mamba, we aim to utilize this feature for conditional prediction. The Gamba block is composed of a Mamba block, two linear projections, and straightforward Prepend and Drop operations:

$$
\begin{aligned}
\mathbf{H}_n &= M_n(\text{Prepend}(\mathbf{P}_c^n \mathbf{T}, \mathbf{P}_x^n \mathbf{X}), \mathbf{G}_{n-1}), \\
\mathbf{G}_n &= \text{Drop}(\mathbf{H}_n, \text{Index}(\mathbf{P}_c^n \mathbf{T}, \mathbf{P}_x^n \mathbf{X})),
\end{aligned}
\tag{2}
$$

where $M_n$ represents the $n$-th vanilla Mamba block. The $\mathbf{P}_c^n \in \mathbb{R}^{D \times C}$ and $\mathbf{P}_x^n \in \mathbb{R}^{D \times C}$ are learnable linear projections for camera and image tokens, respectively, in the $n$-th layer. The operation Prepend involves adding projected camera embeddings and image tokens to the beginning of the sequence before processing through the hidden 3DGS features $\mathbf{G}_{n-1}$ in each layer. Conversely, Drop removes the earlier prepended tokens from the output $\mathbf{H}_n$, based on their indexed positions.

**Gaussian Decoder.** With $N$ stacked Gamba blocks, our GambaFormer adeptly extracts hidden features for each 3DGS token condition on the reference image. A sophisticated Gaussian Decoder then decodes the attributes of each Gaussian $\mathcal{G}_j$. Initially, the output $\mathbf{G}_{N-1}$ from the GambaFormer is input into a shallow MLP, encoding the 3DGS tokens as $\mathbf{Z} = \Phi_\theta$, where $\Phi_\theta$ represents the MLP with learnable parameters $\theta$, and $\mathbf{Z} \in \mathbb{R}^{L \times D}$. Subsequently, separate linear projections are applied to predict each attribute. To precisely predict the $N$ central positions $\mu_j$, we discretize the coordinate space $\mu_j \in [-0.5, 0.5]^3$ into 21 uniformly spaced coordinate points $c_j$. A linear projection $\mathbf{W} \in \mathbb{R}^{21 \times D}$ then maps $\mathbf{Z}$ to $\mathbf{Q} \in \mathbb{R}^{N \times 21}$, and the predicted coordinate can be denoted as:

$$
P(Q_{ij}) = \frac{e^{Q_{ij}}}{\sum_{k=1}^{21} e^{Q_{ik}}}, \quad y_i = \sum_{j=1}^{21} P(Q_{ij}) \times c_j, \quad \text{for } c_j \in \{-0.5, -0.45, \dots, 0.45, 0.5\} \tag{3}
$$

Here, for simplicity, the formulation only considers one axis. $P(Q_{ij})$ represents the softmax probability of the $j$-th coordinate point for the $i$-th token, and $y_i$ denotes the predicted coordinate for the $i$-th token. Opacity $\alpha_j$ is derived using a linear projection followed by a Sigmoid activation. Scale $r_j$ is predicted using a linear projection with a Softplus activation without constrain. Additionally, only the 0-th order of spherical harmonics $c_i$, constrained within the RGB space, is predicted.

### 3.3 Robust Amortized Training.

**Gaussian Parameter Constraints.** Learning accurate 3D positions of Gaussians from a single image presents significant challenges due to the limited geometric information available from a single viewpoint. Prior works [60, 68] has employed point clouds, sampled from ground-truth meshes, as supervision during the initial training phase to prevent model collapse. This form of supervision, however, impedes both the end-to-end training and the scalability of large reconstruction models.

To address these issues, we introduce a *radial mask constraint*, inspired by the notion that images from multiple viewpoints can depict the occupied 3D space of an object. As illustrated in Figure 3, the view mask $\mathbf{x}_{mask}$ is first discretized into a distance field for fast approximation. Using 2D ray casting from the image center to the mask contour, we obtain a set of radial contour distance fields $\mathbf{v} \in \mathbb{R}^U$, where $U$ denotes the number of rays. If the projected 2D center of a Gaussian falls outside these contours, an explicit loss is applied to correct its position:
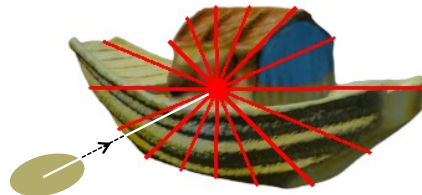


Figure 3: **Radial polygon mask**. Object masks are divided into polygon masks by 2D ray casting from the image center to the contours.

$$
\mathcal{L}_{\text{rdist}} = \mathcal{L}_{\text{MSE}}\left(\text{Interp}(v^k, v^{k+1}), d_j\right), \quad \text{where } d_j = \sqrt{(x_j - c_x)^2 + (y_j - c_y)^2} \tag{4}
$$

Here, $d_j$ represents the distance from the 2D projected center of Gaussian $G_j$ to the image center $(c_x, c_y)$, and $\mathrm{Interp}(v^k, v^{k+1})$ denotes the linearly interpolated ground truth distance between two adjacent rays $v^k$ and $v^{k+1}$, based on the radial angle of the projected 2D Gaussians.

This image-based constraint allows us to eliminate the need for explicit 3D supervision previously required in single-view reconstruction models. By imposing such constraints on the projected 2D Gaussians, the positions of the Gaussians can rapidly converge to rough 3D shapes.

**Training Objective.** Utilizing the efficient tiled rasterizer for 3D Gaussians [20], Gamba is trained end-to-end with image-space reconstruction loss across both reference and novel views. The training loss function is formulated as follows:

$$\mathcal{L}_{train} = \frac{1}{V+1} \sum_{i=0}^{V} \mathcal{L}_{\mathrm{MSE}}(\hat{v}_i^{\mathrm{rgb}}, v_i^{\mathrm{rgb}}) + \lambda_{\mathrm{mask}} \mathcal{L}_{\mathrm{MSE}}(\hat{v}_i^{\alpha}, v_i^{\alpha})$$
$$+ \lambda_{\mathrm{LPIPS}} \mathcal{L}_{\mathrm{LPIPS}}(\hat{v}_i^{\mathrm{rgb}}, v_i^{\mathrm{rgb}}) + \lambda_{\mathrm{rdist}} \mathcal{L}_{\mathrm{rdist}}, \tag{5}$$

where $\hat{v}_i^{\mathrm{rgb}}$ and $\hat{v}_i^{\alpha}$ denote the predicted RGB image and alpha mask, respectively, rendered from the predicted 3D Gaussians $\{\mathcal{G}\}$, while $v_i^{\mathrm{rgb}}$ and $v_i^{\alpha}$ represent the corresponding ground truth. $\mathcal{L}_{\mathrm{LPIPS}}$ encompasses the VGG-based perceptual loss for image fidelity [65], and $\mathcal{L}_{\mathrm{rdist}}$ is the radial mask constraint loss defined in Eq. (4). Additionally, the $\lambda_{mask}$, $\lambda_{\mathrm{LPIPS}}$ and $\lambda_{\mathrm{rdist}}$ are balancing factors. These losses are applied across $V + 1$ viewpoints, including the input reference viewpoint and $V$ novel viewpoints to supervise geometry and texture traininig together.

Inspired by the coarse-to-fine optimization strategies employed in SDS-based image-to-3D reconstructions [37, 63, 38], we progressively increase the number of views from 2 to 6 during training. This strategy not only reduces the computational load but also enhances the model's robustness by allowing gradual adaptation from geometry to detailed textures.

## 4 Experiments

### 4.1 Implementation Details

**Datasets.** Following previous work [26, 68], we utilized a filtered LVIS subset of the Objaverse dataset [5] for pre-training Gamba, which comprises around 40k 3D models across 1,156 categories. We filtered this subset by intersecting it with 3D objects rendered in the G-buffer Objaverse [39], resulting in a final training set of approximately 20k high-quality 3D objects. Additionally, to improve model generalization, the camera poses for rendered objects during training are normalized to a unit distance [47]. For evaluation, our trained model was qualitatively assessed using web images. Quantitatively, we conducted comparisons on the Google Scanned Object (GSO) dataset [6], selecting totally 60 3D objects and rendering a single view of each at a spatial resolution of $512 \times 512$ for comprehensive evaluation.

**Network Architecture.** The GambaFormer architecture comprises $N = 14$ Gamba blocks, each with hidden dimensions $D = 512$. For the condition image tokens, we employ the pretrained DINO v2 model [35] as our image tokenizer, which extracts semantic feature tokens of length $K = 576$ from the reference image. To construct the 3DGS token sequence, the reference image, with a spatial resolution of $512 \times 512$, is initially processed with a convolution kernel with $p = 8$. This step tokenizes the image into 4096 tokens. Subsequently, four pre-defined scan orders [28] are applied sequentially to expand this into a token sequence of length $L = 16384$. The 3DGS embeddings **E** are learnable positional embeddings and correspond to 16384 3D Gaussians, matching the length of the above tokens. The Gaussian Decoder employs a straightforward MLP architecture with a single hidden layer, complemented by separate linear projections for each attribute of the Gaussians.

**Pre-training.** Gamba is trained on 16 NVIDIA A100 (80G) GPUs with a batch size of 256 over approximately 40 hours for totally 400 epochs. We employ the AdamW optimizer [29] with an initial learning rate of 1e-3 and a weight decay of 0.05. Gradient clipping of value 1.0 was implemented to maintain the $L_2$ norm of gradients. The loss weight settings are as follows: $\lambda_{\mathrm{mask}} = 1$ for the mask loss and $\lambda_{\mathrm{LPIPS}} = 0.5$ for LPIPS [65] loss. Additionally, the loss weight $\lambda_{\mathrm{rdist}}$ for the radial mask constraints is initially set at 0.1 and is gradually decayed to 0 over the first 10 epochs, aiming to ensure that the predicted positions of 3D Gaussians converge within a reasonable range in 3D space.

Figure 4: Qualitative Comparison with large reconstruction models.

**Inference.** During inference, Gamba only takes an arbitary RGB image as input, where the foreground object is segmented using a pre-trained segmentation model [21] and subsequently recentered. Gamba employs a default camera pose with both zero elevation and azimuth, which is used to produce camera tokens and Plücker ray inputs. Gamba efficiently predicts $16,384$ Gaussians for a single 3D object in a feed-forward manner. Remarkably, this process requires only about *8 GB* of GPU memory and completes in less than *0.05 second* on a single NVIDIA A100 (80G) GPU, making it well-suited for online deployment scenarios.

## 4.2 Experimental Protocol

**Baseline.** We benchmark Gamba against previous single-view reconstruction methods, particularly those in the stream of large reconstruction models. These models are typically trained on a large number of rendered multi-view images and are designed for efficient feed-forward inference. The first work, LRM [18] utilizes a transformer-based architecture to predict a Tri-Plane representation from a single image. Triplane-Meets-Gaussian (TGS) combines the fast rendering capabilities of 3DGS by initially predicting a point cloud of $16,384$ for a 3D object, followed by predicting other Gaussian attributes using another transformer network. In contrast, another stream is SDS-based methods with iterative optimization like DreamGaussian [48]. These approaches leverage a multi-view diffusion model, Zero-1-to-3 [27], trained on the Objaverse-XL dataset [4], to produce multi-view images conditioned on a single image and relative camera poses. Another notable method, One-2-3-45, bypasses the need for costly optimization by utilizing a generalizable SparseNeuS [51, 54] to directly predict Signed Distance Functions (SDF) from generated multi-view images.

**Qualitative Comparisons.** Figures. 4 and 5 demonstrate Gamba's capability to maintain reasonable geometry and plausible textures in reconstructing various 3D objects. In contrast, reconstructions by most baseline methods suffer from multi-view inconsistency and geometric distortion. Despite LRM being trained on a dataset that is $50\times$ larger than ours, it frequently exhibits warped geometries (row 1, 2, 5 in Figure 4). Compared to TGS, which also employs 3DGS for representation, Gamba consistently delivers better texture reconstruction (especially row 1, 3, 4 in Fig. 4). Fig. 5 further highlights that, while both One-2-3-45 and DreamGaussian leverage the advanced Zero-1-to-3-XL model [27, 4], their reconstruction still exhibit artifacts with multi-view inconsistency and geometric

7

distortion (rows 1 and 3) alongside blurred textures (rows 2, 3, and 4). This comparative analysis underscores Gamba's robustness and superior performance in single-view 3D reconstruction.



Figure 5: Comparison with Zero-1-to-3 [27] based single-view 3D reconstruction methods, including feed-forward only method One-2-3-45 [26] and optimization-based DreamGaussian [48].

**Quantitative Comparisons.** Following prior works, we adopt PSNR, SSIM, and LPIPS metrics to evaluate the quality of novel view synthesis. Additionally, we employ Chamfer Distance and Volume IoU between ground truth and reconstruction to quality of geometry reconstruction. Gamba is compared against 2 SDS-based methods: Zero-1-to-3 [27] and DreamGaussian [48], as well as 5 feed-forward only methods: One-2-3-45 [26], the point cloud diffusion model Point-E [34], the NerF parameter diffusion model Shap-E [19], the pioneering large reconstruction model LRM [18], and TGS which integrates LRM with 3DGS [68]. The results, presented in Table 1, indicate that Gamba is competitive in both texture and geometry reconstruction.

Table 1: **Quantitative results.** We evaluate novel view synthesis in terms of PSNR↑/ SSIM↑/LPIPS↓ and geometry reconstruction in terms of Chamfer Distance↓/Volume IoU↑.

| | PSNR ↑ | SSIM ↑ | LPIPS ↓ | Chamfer Dist. ↓ | Volume IoU ↑ | Time ↓ |
|---|---|---|---|---|---|---|
| One-2-3-45 [26] | 14.93 | 0.72 | 0.28 | 0.0683 | 0.3485 | 40s |
| Point-E [34] | - | - | - | 0.0515 | 0.2643 | 78s |
| Shap-E [19] | - | - | - | 0.0579 | 0.3228 | 27s |
| Zero-1-to-3 [27] | 18.28 | 0.76 | 0.19 | 0.0385 | 0.3786 | 1800s |
| DreamGaussian [48] | 21.65 | 0.82 | 0.14 | 0.0341 | 0.3615 | 70 s |
| LRM [18] | 21.35 | 0.82 | 0.15 | 0.0325 | 0.3872 | 0.5s |
| TGS [68] | 22.68 | 0.85 | 0.12 | 0.0257 | 0.4121 | 0.2s |
| **Ours** | **24.74** | **0.91** | **0.08** | **0.0232** | **0.4289** | **0.05s** |

**Inference Runtime** We showcase the inference runtime required to generate a 3D asset in Table 1, where the timing is recorded using the default hyper-parameters for each method on a single NVIDIA A100 GPU (80G). Remarkably, our Gamba outperforms optimization-based approaches like Zero-1-to-3 [27] in terms of speed, being several orders of magnitude faster than those optimization-based methods and surpass other feed-forward models as well, thanks to the efficient backbone design.

# 5 Ablation and Discussion

In ablation studies, all experiments is conducted on a randomly selected subset from G-buffer [39] Objaverse, around totally 10k training data, and 500 objects in this subset are left for evaluation. And all models are trained 100 epochs only for evaluation.

**Q1:** *What impacts performance of Gamba in terms of component-wise contributions?* We discarded each core component of Gamba to validate its component-wise effectiveness. The results are described in Table 2 by comparing their predicted views with ground-truth multi-view image set on the evaluation set.

**A1:** In an end-to-end, multi-component pipeline, we observed that the exclusion of any component from Gamba resulted in a significant degradation in performance. Specifically, we first remove the loss term of the radial mask constraint by setting $\lambda_{\text{rdist}} = 0$ in Eq. (5) for the "w/o radial mask constraint". We find that the training is prone to collapse, i.e., the 3D position of predicted Gaussians confines to a small sphere, after which all predicted opacities $\alpha$ become 0, leading to Gaussians becoming invisible in 2D rendering outputs. This removal thus produces a catastrophic performance drop; the evaluation PSNR is only 12.72. Furthermore, in comparison to prior works such as AGG [60] and Triplane-Meets-Gaussian [68], which exclusively utilized learnable positional embeddings as 3D tokens, our model was also evaluated under this variant. The "w/o additive 3DGS tokens" scenario means taking the 3DGS tokens solely from learnable embeddings, i.e., setting $\mathbf{G} = \mathbf{E}$ in Eq. (1); the evaluation PSNR degraded to 20.35 in quantitative evaluation. In qualitative comparison, we find that the reconstruction after this modification tends to produce over-smoothed texture, and the generalization of the model to different objects is also degraded. Finally, we removed the prepending operation of the conditional camera tokens and image tokens. The evaluation results, shown in "w/o Prepending," reveal that removing this component leads to a 1.24 dB drop in terms of PSNR, which demonstrates the necessity of this prepending operation.

Table 2: Ablation Studies of component-wise contribution.

| Model variants | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|
| w/o radial mask constraint | 12.72 | 0.58 | 0.47 |
| w/o additive 3DGS tokens | 20.35 | 0.79 | 0.16 |
| w/o Prepending | 22.57 | 0.85 | 0.10 |
| Full Model | 23.81 | 0.88 | 0.12 |

**Q2:** *Why don't we apply more Gaussians?* Mamba [10] exhibits linear computational complexity and memory consumption with respect to token length. We empirically validated that $L = 16384$ Gaussians are sufficient for amortized 3D reconstruction in our end-to-end training framework.

**A2:** We illustrate the memory consumption of Mamba and Transformer in Figure 7. Existing 3DGS-based amortized 3D reconstruction models [60, 68], which are all built on transformer architectures, have a maximum of 4096 3D Gaussians due to the $O(N^2)$ memory consumption of transformers. These models employ local attention or 1D-CNN to predict offsets over the 4096 Gaussians to achieve a higher resolution, such as 16384 Gaussians from 4096 base Gaussians. However, these models are trained in two stages, which limits the performance of the higher resolution Gaussians by the capabilities of the first-stage model. In contrast, our method, which is trained end-to-end, directly applies 16384 Gaussians in a single forward. Additionally, we quantitatively reconstructed 500 objects using original 3D Gaussian Splatting with multi-view images and statistically analyzed their total Gaussians after iterative reconstruction. We found that 87% of the 3D objects maintained Gaussian counts between 10000 and 20000. Thus, using $L = 16384$ for Gaussian numbers is sufficient in our Gamba framework.

**Q3:** *Why do we not use Tri-plane representation?* Almost all previous amortized single view reconstruction models integrate Tri-plane representation to encode textures or the geometry of 3D objects. We also explored adding the Tri-plane representation to our Gamba model for ablation study.

**A3:** In Figure 6 (a), we show that our Tri-Plane meets-Gaussian structure. A Siamese GambaFormer architecture is adopted where the left GambaFormer predicts positions $\mu_j$, opacity $\alpha_j$, and scales $r_j$. Similarly, the right branch adopts a methodology as described in Triplane-Meets-Gaussian [68]. Here, the predicted positions of Gaussians $\mu_j$ serve as queries to extract hidden features from the predicted hidden Tri-Plane features. Subsequently, these queried features from three planes are concatenated to decode the color $c_j$ of each Gaussian. Ultimately, the prediction results from these two

branches are combined and rendered into multi-view images for supervision, employing the same supervision as in Eq. (5) during training. The evaluation of PSNR on the selected G-buffer Objaverse subset is shown in Figure 6 (b). Despite having doubled parameters, the PSNR of "Triplane-Gamba" is significantly worse than our base version, and its PSNR sharply degrades at the 40th epoch. This ablation study demonstrates that the Tri-Plane representation is not necessary in our Gamba Framework.
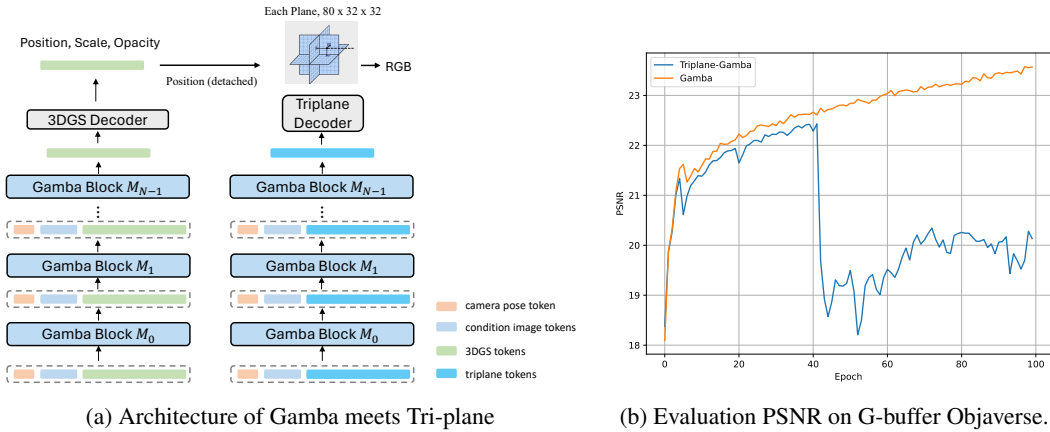


(a) Architecture of Gamba meets Tri-plane

(b) Evaluation PSNR on G-buffer Objaverse.

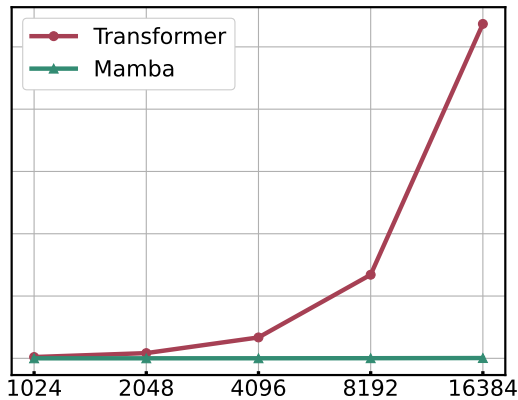Figure 6: Ablation study of Tri-plane meets Gamba.



Figure 7: Mamba vs Transformer memory consumption comparison over token length.

**Q4:** *Are there any other alternatives to construct the $L$ 3DGS tokens?* In our initial experiment (Q1), we explored constructing the $L$ 3DGS tokens solely using learnable positional embeddings $\mathbf{E}$. This approach resulted in a significant performance drop, indicating the insufficiency of relying purely on embeddings for constructing 3DGS tokens. In the current methodology, as detailed in Eq. (1), we expand the tokenized images by scanning them four times to produce $L = 16384$ tokens. This method, while simple yet effective, raises concerns that the reconstructed 3D objects might merely duplicate the conditional image. To verify this concerns, we explore to construct the $L$ 3DGS tokens from a mixed strategy.

**A4**: Specifically, we initially construct $1024$ tokens using a convolution with a kernel size of $p = 8$. These tokens are then expanded to $4096$ through four scans, as specified in Eq. (1). The remaining $12,288$ tokens are constructed directly from pure positional embeddings $\mathbf{E}$. This method's qualitative results are illustrated in Figure 8, where the rendered 2D views from the predicted 3D Gaussians exhibit noticeable blurring across all views, including the reference view. This degradation in quality relative to our baseline methods indicates that relying solely on positional embeddings can lead to significant artifacts, affecting both geometry and texture clarity. Our analysis suggests that the difficulty in achieving convergence with positional embeddings during training, combined with their limited capacity to encapsulate relevant information from the conditional image. Conversely, em-

10

bedding the image directly onto each 3DGS token through four different scan order helps the model progressively predict finer details of the 3D object, aligning more closely with the unidirectional scan on the input 3DGS tokens.



Input                                    Predicted views

Figure 8: Qualitative evaluation of mixed 3D GS tokens. Predicted 3D objects on 2D views rendering appears to be blurred even at the reference viewpoint.

## 6   Conclusion

In this work, we present Gamba, the first end-to-end trained, amortized 3D reconstruction model from single-view image. Our proposed Gamba, different from previous methods reliant on SDS and NeRF, marries 3D Gaussian splatting and Mamba to address the challenges of high memory requirements and heavy rendering process. Our key insight is the relationship between the 3DGS generation process and the sequential mechanism of Mamba. Additionally, Gamba integrates several techniques for training stability. Through extensive qualitative comparisons and quantitative evaluations, we show that our Gamba is promising and competitive with several orders of magnitude speedup in single-view 3D reconstruction.

# Gamba: Marry Gaussian Splatting with Mamba for Single-View 3D Reconstruction

## Supplementary Material

## 7 Preliminary of State Space Models

Utilizing ideas from the control theory [7], the integration of linear state space equations with deep learning has been widely employed to tackle the modeling of sequential data. The promising property of linearly scaling with sequence length in long-range dependency modeling has attracted great interest from searchers. Pioneered by LSSL [13] and S4 [12], which utilize linear state space equations for sequence data modeling, follow-up works mainly focus on memory efficiency [12], fast training speed [14, 11] and better performance [31, 50]. More recently, Mamba [10] integrates a selective mechanism and efficient hardware design, outperforms Transformers [49] on natural language and enjoys linear scaling with input length. Building on the success of Mamba, Vision Mamba [67] and VMamba [28] leverage the bidirectional Vim Block and the Cross-Scan Module respectively to gain data-dependent global visual context for visual representation; U-Mamba [30] and Vm-unet [41] further bring Mamba into the field of medical image segmentation. PointMamba [24] and Point Cloud Mamba [66] adapt Mamba for point cloud understanding through reordering and serialization strategy.

**State Space Models (SSMs)** [12] have emerged as a powerful tool for modeling and analyzing complex physical systems, particularly those that exhibit linear time-invariant (LTI) behavior. The core idea behind SSMs is to represent a system using a set of first-order differential equations that capture the dynamics of the system's state variables. This representation allows for a concise and intuitive description of the system's behavior, making SSMs well-suited for a wide range of applications. The general form of an SSM can be expressed as follows:

$$
\begin{aligned}
\dot{h}(t) &= Ah(t) + Bx(t), \\
y(t) &= Ch(t) + Dx(t).
\end{aligned}
\tag{6}
$$

where $h(t)$ denotes the state vector of the system at time $t$, while $\dot{h}(t)$ denotes its time derivative. The matrices $A$, $B$, $C$, and $D$ encode the relationships between the state vector, the input signal $x(t)$, and the output signal $y(t)$. These matrices play a crucial role in determining the system's response to various inputs and its overall behavior.

One of the challenges in applying SSMs to real-world problems is that they are designed to operate on continuous-time signals, whereas many practical applications involve discrete-time data. To bridge this gap, it is necessary to discretize the SSM, converting it from a continuous-time representation to a discrete-time one. The discretized form of an SSM can be written as:

$$
\begin{aligned}
h_k &= \bar{A}h_{k-1} + \bar{B}x_k, \\
y_k &= \bar{C}h_k + \bar{D}x_k.
\end{aligned}
\tag{7}
$$

Here, $k$ represents the discrete time step, and the matrices $\bar{A}$, $\bar{B}$, $\bar{C}$, and $\bar{D}$ are the discretized counterparts of their continuous-time equivalents. The discretization process involves sampling the continuous-time input signal $x(t)$ at regular intervals, with a sampling period of $\Delta$. This leads to the following relationships between the continuous-time and discrete-time matrices:

$$
\begin{aligned}
\bar{A} &= (I - \Delta/2 \cdot A)^{-1}(I + \Delta/2 \cdot A), \\
\bar{B} &= (I - \Delta/2 \cdot A)^{-1}\Delta B, \\
\bar{C} &= C.
\end{aligned}
\tag{8}
$$

**Selective State Space Models** [9] are proposed to address the limitations of traditional SSMs in adapting to varying input sequences and capturing complex, input-dependent dynamics. The key innovation in Selective SSMs is the introduction of a selection mechanism that allows the model to efficiently select data in an input-dependent manner, enabling it to focus on relevant information and ignore irrelevant inputs. The selection mechanism is implemented by parameterizing the SSM matrices $\bar{B}$, $\bar{C}$, and $\Delta$ based on the input $x_k$. This allows the model to dynamically adjust its behavior

depending on the input sequence, effectively filtering out irrelevant information and remembering relevant information indefinitely.

## 8  Limitations

While Gamba achieves remarkable speed and promising results in 3D reconstruction, there are still some limitations. Firstly, the reconstruction quality is highly dependent on the input image; if the input lacks sufficient geometric information, Gamba struggles to produce accurate reconstructions. Secondly, the texture of the reconstructed back view is often smoothed, as it may significantly differ from what is visible in the input image. This issue stems from Gamba being trained under a prediction-only paradigm rather than as a generative model, which limits its ability to infer unseen views of 3D objects. In future, we might explore incorporating a diffusion training to enhance Gamba's generative capabilities. Thirdly, the scalability of the Gamba model remains underexplored. Currently, the model utilizes only 14 blocks, with its parameters amounting to merely about 10% of those employed in existing large reconstruction models. We envision this work as paving the way for alternative approaches in amortized 3D reconstruction, inviting further exploration into scalable architectures.

# References

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

[2] David Charatan, Sizhe Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction. *arXiv preprint arXiv:2312.12337*, 2023.

[3] Guikun Chen and Wenguan Wang. A survey on 3d gaussian splatting. *arXiv preprint arXiv:2401.03890*, 2024.

[4] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, Eli VanderBilt, Aniruddha Kembhavi, Carl Vondrick, Georgia Gkioxari, Kiana Ehsani, Ludwig Schmidt, and Ali Farhadi. Objaverse-xl: A universe of 10m+ 3d objects. *arXiv preprint arXiv:2307.05663*, 2023.

[5] 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*, pages 13142–13153, 2023.

[6] 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)*, pages 2553–2560. IEEE, 2022.

[7] William Glasser. *Control theory*. Harper and Row New York, 1985.

[8] Alexandros Graikos, Nikolay Malkin, Nebojsa Jojic, and Dimitris Samaras. Diffusion models as plug-and-play priors. *Advances in Neural Information Processing Systems*, 35:14715–14728, 2022.

[9] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.

[10] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.

[11] Albert Gu, Karan Goel, Ankit Gupta, and Christopher Ré. On the parameterization and initialization of diagonal state space models. *Advances in Neural Information Processing Systems*, 35:35971–35983, 2022.

[12] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021.

[13] Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, and Christopher Ré. Combining recurrent, convolutional, and continuous-time models with linear state space layers. *Advances in neural information processing systems*, 34:572–585, 2021.

[14] Albert Gu, Isys Johnson, Aman Timalsina, Atri Rudra, and Christopher Ré. How to train your hippo: State space models with generalized orthogonal basis projections. *arXiv preprint arXiv:2206.12037*, 2022.

[15] Faiza Gul, Wan Rahiman, and Syed Sahal Nazli Alhady. A comprehensive study for robot navigation techniques. *Cogent Engineering*, 6(1):1632046, 2019.

[16] Yuan-Chen Guo, Ying-Tian Liu, Ruizhi Shao, Christian Laforte, Vikram Voleti, Guan Luo, Chia-Hao Chen, Zi-Xin Zou, Chen Wang, Yan-Pei Cao, and Song-Hai Zhang. threestudio: A unified framework for 3d content generation. https://github.com/threestudio-project/threestudio, 2023.

[17] Susung Hong, Donghoon Ahn, and Seungryong Kim. Debiasing scores and prompts of 2d diffusion for robust text-to-3d generation. *arXiv preprint arXiv:2303.15413*, 2023.

[18] 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. *arXiv preprint arXiv:2311.04400*, 2023.

[19] Heewoo Jun and Alex Nichol. Shap-e: Generating conditional 3d implicit functions. *arXiv preprint arXiv:2305.02463*, 2023.

[20] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (ToG)*, 42(4):1–14, 2023.

[21] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023.

[22] 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. *arXiv preprint arXiv:2311.06214*, 2023.

[23] Dingkang Liang, Xin Zhou, Xinyu Wang, Xingkui Zhu, Wei Xu, Zhikang Zou, Xiaoqing Ye, and Xiang Bai. Pointmamba: A simple state space model for point cloud analysis. *arXiv preprint arXiv:2402.10739*, 2024.

[24] Dingkang Liang, Xin Zhou, Xinyu Wang, Xingkui Zhu, Wei Xu, Zhikang Zou, Xiaoqing Ye, and Xiang Bai. Pointmamba: A simple state space model for point cloud analysis. *arXiv preprint arXiv:2402.10739*, 2024.

[25] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 300–309, 2023.

[26] Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Zexiang Xu, Hao Su, et al. One-2-3-45: Any single image to 3d mesh in 45 seconds without per-shape optimization. *arXiv preprint arXiv:2306.16928*, 2023.

[27] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9298–9309, 2023.

[28] Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, and Yunfan Liu. Vmamba: Visual state space model. *arXiv preprint arXiv:2401.10166*, 2024.

[29] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

[30] Jun Ma, Feifei Li, and Bo Wang. U-mamba: Enhancing long-range dependency for biomedical image segmentation. *arXiv preprint arXiv:2401.04722*, 2024.

[31] Harsh Mehta, Ankit Gupta, Ashok Cutkosky, and Behnam Neyshabur. Long range language modeling via gated state spaces. *arXiv preprint arXiv:2206.13947*, 2022.

[32] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.

[33] Charlie Nash, Yaroslav Ganin, SM Ali Eslami, and Peter Battaglia. Polygen: An autoregressive generative model of 3d meshes. In *International conference on machine learning*, pages 7220–7229. PMLR, 2020.

[34] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. Point-e: A system for generating 3d point clouds from complex prompts. *arXiv preprint arXiv:2212.08751*, 2022.

[35] Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2023.

[36] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.

[37] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022.

[38] Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren, Aliaksandr Siarohin, Bing Li, Hsin-Ying Lee, Ivan Skorokhodov, Peter Wonka, Sergey Tulyakov, et al. Magic123: One image to high-quality 3d object generation using both 2d and 3d diffusion priors. *arXiv preprint arXiv:2306.17843*, 2023.

[39] Lingteng Qiu, Guanying Chen, Xiaodong Gu, Qi zuo, Mutian Xu, Yushuang Wu, Weihao Yuan, Zilong Dong, Liefeng Bo, and Xiaoguang Han. Richdreamer: A generalizable normal-depth diffusion model for detail richness in text-to-3d. *arXiv preprint arXiv:2311.16918*, 2023.

[40] 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*, pages 10684–10695, 2022.

[41] Jiacheng Ruan and Suncheng Xiang. Vm-unet: Vision mamba unet for medical image segmentation. *arXiv preprint arXiv:2402.02491*, 2024.

[42] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view diffusion for 3d generation. *arXiv preprint arXiv:2308.16512*, 2023.

[43] J. Ryan Shue, Eric Ryan Chan, Ryan Po, Zachary Ankner, Jiajun Wu, and Gordon Wetzstein. 3d neural field generation using triplane diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20875–20886, June 2023.

[44] Tianyu Sun, Guodong Zhang, Wenming Yang, Jing-Hao Xue, and Guijin Wang. Trosd: A new rgb-d dataset for transparent and reflective object segmentation in practice. *IEEE Transactions on Circuits and Systems for Video Technology*, 2023.

[45] Stanislaw Szymanowicz, Christian Rupprecht, and Andrea Vedaldi. Splatter image: Ultra-fast single-view 3d reconstruction. *arXiv preprint arXiv:2312.13150*, 2023.

[46] Jiaxiang Tang. Stable-dreamfusion: Text-to-3d with stable-diffusion, 2022. https://github.com/ashawkey/stable-dreamfusion.

[47] 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. *arXiv preprint arXiv:2402.05054*, 2024.

[48] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023.

[49] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[50] Jue Wang, Wentao Zhu, Pichao Wang, Xiang Yu, Linda Liu, Mohamed Omar, and Raffay Hamid. Selective structured state-spaces for long-form video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6387–6397, 2023.

[51] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689*, 2021.

[52] Peng Wang and Yichun Shi. Imagedream: Image-prompt multi-view diffusion for 3d generation. *arXiv preprint arXiv:2312.02201*, 2023.

[53] Peng Wang, Hao Tan, Sai Bi, Yinghao Xu, Fujun Luan, Kalyan Sunkavalli, Wenping Wang, Zexiang Xu, and Kai Zhang. Pf-lrm: Pose-free large reconstruction model for joint pose and shape prediction. *arXiv preprint arXiv:2311.12024*, 2023.

[54] Yiming Wang, Qin Han, Marc Habermann, Kostas Daniilidis, Christian Theobalt, and Lingjie Liu. Neus2: Fast learning of neural implicit surfaces for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3295–3306, 2023.

[55] Yuxuan Wang, Xuanyu Yi, Zike Wu, Na Zhao, Long Chen, and Hanwang Zhang. View-consistent 3d editing with gaussian splatting. *arXiv preprint arXiv:2403.11868*, 2024.

[56] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolific-dreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *arXiv preprint arXiv:2305.16213*, 2023.

[57] Haohan Weng, Tianyu Yang, Jianan Wang, Yu Li, Tong Zhang, CL Chen, and Lei Zhang. Consistent123: Improve consistency for one image to 3d object synthesis. *arXiv preprint arXiv:2310.08092*, 2023.

[58] Tong Wu, Jiarui Zhang, Xiao Fu, Yuxin Wang, Jiawei Ren, Liang Pan, Wayne Wu, Lei Yang, Jiaqi Wang, Chen Qian, et al. Omniobject3d: Large-vocabulary 3d object dataset for realistic perception, reconstruction and generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 803–814, 2023.

[59] Zike Wu, Pan Zhou, Xuanyu Yi, Xiaoding Yuan, and Hanwang Zhang. Consistent3d: Towards consistent high-fidelity text-to-3d generation with deterministic sampling prior. *arXiv preprint arXiv:2401.09050*, 2024.

[60] Dejia Xu, Ye Yuan, Morteza Mardani, Sifei Liu, Jiaming Song, Zhangyang Wang, and Arash Vahdat. Agg: Amortized generative 3d gaussians for single image to 3d. *arXiv preprint arXiv:2401.04099*, 2024.

[61] Yinghao Xu, Hao Tan, Fujun Luan, Sai Bi, Peng Wang, Jiahao Li, Zifan Shi, Kalyan Sunkavalli, Gordon Wetzstein, Zexiang Xu, et al. Dmv3d: Denoising multi-view diffusion using 3d large reconstruction model. *arXiv preprint arXiv:2311.09217*, 2023.

[62] Xuanyu Yi, Jiajun Deng, Qianru Sun, Xian-Sheng Hua, Joo-Hwee Lim, and Hanwang Zhang. Invariant training 2d-3d joint hard samples for few-shot point cloud recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14463–14474, 2023.

[63] Xuanyu Yi, Zike Wu, Qingshan Xu, Pan Zhou, Joo-Hwee Lim, and Hanwang Zhang. Diffusion time-step curriculum for one image to 3d generation. *arXiv preprint arXiv:2404.04562*, 2024.

[64] Xianggang Yu, Mutian Xu, Yidan Zhang, Haolin Liu, Chongjie Ye, Yushuang Wu, Zizheng Yan, Chenming Zhu, Zhangyang Xiong, Tianyou Liang, et al. Mvimgnet: A large-scale dataset of multi-view images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9150–9161, 2023.

[65] 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*, pages 586–595, 2018.

[66] Tao Zhang, Xiangtai Li, Haobo Yuan, Shunping Ji, and Shuicheng Yan. Point could mamba: Point cloud learning via state space model, 2024.

[67] Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. Vision mamba: Efficient visual representation learning with bidirectional state space model. *arXiv preprint arXiv:2401.09417*, 2024.

[68] 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. *arXiv preprint arXiv:2312.09147*, 2023.

[69] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Ewa splatting. *IEEE Transactions on Visualization and Computer Graphics*, 8(3):223–238, 2002.