

L3DG: Latent 3D Gaussian Diffusion

BARBARA ROESSLE, Technical University of Munich, Germany

NORMAN MÜLLER, Meta Reality Labs Zurich, Switzerland

LORENZO PORZI, Meta Reality Labs Zurich, Switzerland

SAMUEL ROTA BULÒ, Meta Reality Labs Zurich, Switzerland

PETER KONTSCHIEDER, Meta Reality Labs Zurich, Switzerland

ANGELA DAI, Technical University of Munich, Germany

MATTHIAS NIESSNER, Technical University of Munich, Germany



Fig. 1. L3DG learns a compressed latent space of 3D Gaussian representations and efficiently synthesizes novel scenes via diffusion in latent space. This approach makes L3DG scalable to room-size scenes, which are generated from pure noise leading to geometrically realistic scenes of 3D Gaussians that can be rendered in real-time. Above results are from our model trained on 3D-FRONT; we visualize the 3D Gaussian ellipsoids and show renderings.

We propose L3DG, the first approach for generative 3D modeling of 3D Gaussians through a latent 3D Gaussian diffusion formulation. This enables effective generative 3D modeling, scaling to generation of entire room-scale scenes which can be very efficiently rendered. To enable effective synthesis of 3D Gaussians, we propose a latent diffusion formulation, operating in a compressed latent space of 3D Gaussians. This compressed latent space is learned by a vector-quantized variational autoencoder (VQ-VAE), for which we employ a sparse convolutional architecture to efficiently operate on room-scale scenes. This way, the complexity of the costly generation process via diffusion is substantially reduced, allowing higher detail on object-level generation, as well as scalability to large scenes. By leveraging the 3D Gaussian representation, the generated scenes can be rendered from arbitrary viewpoints in real-time. We demonstrate that our approach significantly improves visual quality over prior work on unconditional object-level radiance field synthesis and showcase its applicability to room-scale scene generation.

CCS Concepts: • Computing methodologies → Rendering; Neural networks.

Additional Key Words and Phrases: Generative 3D scene modeling, 3D gaussian splatting, latent diffusion

1 INTRODUCTION

Generation of 3D content provides the foundation for many computer graphics applications, from asset creation for video games and films to augmented and virtual reality and creating immersive visual media. In recent years, volumetric rendering [Kajiya and Von Herzen 1984; Kerbl et al. 2023; Mildenhall et al. 2020] has

become a powerful scene representation for 3D content, enabling impressive photorealistic rendering, as it yields effective gradient propagation. 3D Gaussians [Kerbl et al. 2023] have become a particularly popular representation for volumetric rendering that leverages the traditional graphics pipeline in order to obtain high-fidelity renderings at real-time rates. This combination of fast rendering speed and smooth gradients through the optimization, makes 3D Gaussians an ideal candidate for generative 3D modeling.

Inspired by the success of generative modeling for neural radiance fields of single objects [Chan et al. 2022; Chen et al. 2023; Müller et al. 2023], we aim to design a generative model for 3D Gaussians, which can provide a more scalable, rendering-efficient representation for 3D generative modeling. Unfortunately, such generative modeling of 3D Gaussians remains challenging. In particular, this requires a joint understanding of both scene structure as well as the intricacies of realistic appearance, for varying-sized scenes. Moreover, 3D Gaussians are irregularly structured sets, typically containing large quantities of varying numbers of Gaussians, which a generative model must unify into an effective latent manifold. This necessitates a flexible, scalable learned feature representation from which a generative model can be trained.

We thus propose a new generative approach for unconditional synthesis of 3D Gaussians, as a representation that enables high-fidelity view synthesis for both small-scale single objects using ~8k Gaussians, and enables effective scaling to room-scale scenes

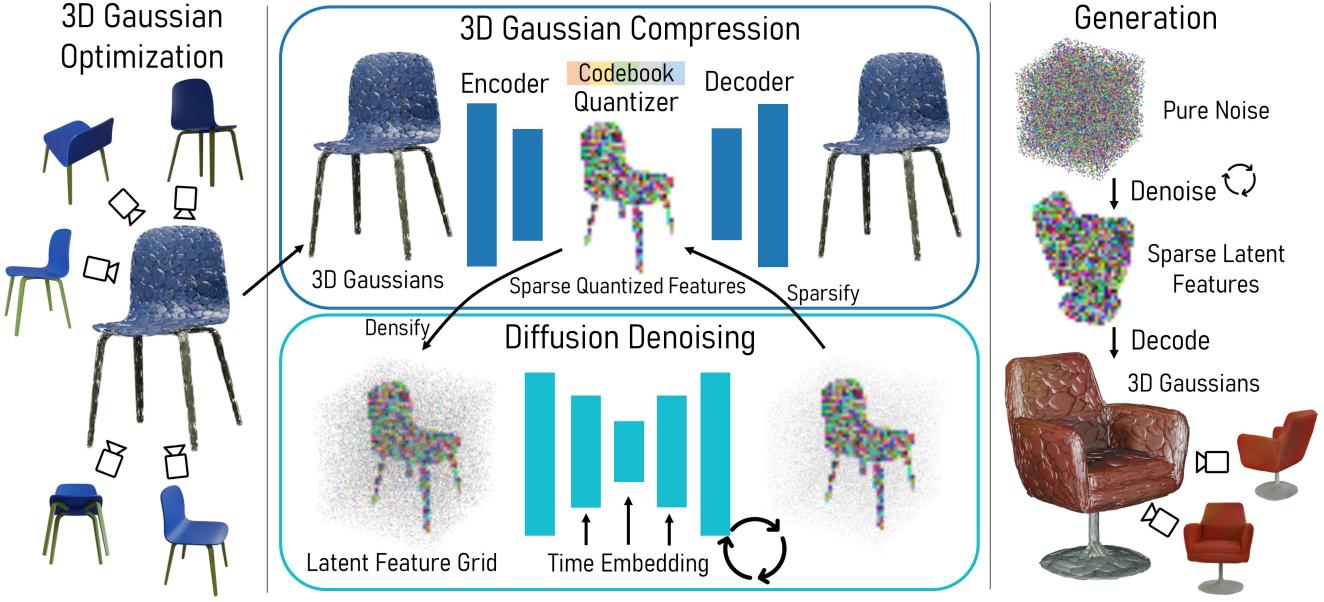


Fig. 2. L3DG method overview: our 3D Gaussian compression model learns to compress 3D Gaussians into sparse quantized features using sparse convolutions and vector-quantization at the bottleneck (VQ-VAE). This allows our 3D diffusion model to efficiently operate on the compressed latent space. At test time, novel scenes are generated by denoising in latent space, which can be sparsified and decoded to high quality 3D Gaussians.

with $\sim 200k$ Gaussians. To facilitate synthesis of large-scale environments, we formulate a latent 3D Gaussian diffusion process. We learn a compressed latent space of 3D Gaussians on a hybrid sparse grid representation for 3D Gaussians, where each sparse voxel encodes a corresponding 3D Gaussian. This latent space is trained as a vector-quantized variational autoencoder (VQ-VAE), and its efficient encoding of 3D Gaussians enables flexible representation scaling from objects to 3D rooms. We then train the generation process through diffusion on this latent 3D Gaussian space, enabling high-fidelity synthesis of 3D Gaussians representing room-scale scenes. Experiments on both object and 3D scene data show that our approach not only produces higher quality synthesis of objects than state of the art, but also much more effectively scales to large scenes, producing 3D scene generation with realistic view synthesis. Our latent 3D Gaussian diffusion improves the FID metric by $\sim 45\%$ compared to DiffRF on PhotoShape [Park et al. 2018].

In summary, our contributions are:

- the first approach to model 3D Gaussians as a generative latent diffusion model, enabling effective synthesis of 3D Gaussian representations of room-scale scenes that yield realistic view synthesis.
- our latent 3D Gaussian diffusion formulation enables flexible generative modeling on a compressed latent space constructed by sparse 3D convolutions, capturing both high-fidelity objects as well as larger, room-scale scenes.

2 RELATED WORK

Our work addresses the problem of unconditional generation of 3D objects and scenes. We review below related works categorizing them based on the type of generative model.

GAN-based. Generative Adversarial Networks [Goodfellow et al. 2014] (GAN) have found successful application in the generation of 3D assets. The generator maps random noise to the target 3D representation, from which images are rendered given camera poses and supervised with the discriminator. Accordingly, methods in this category can be trained without having explicit access to 3D ground-truth, but rely only on posed images. Several 3D representations have been considered in different works ranging from simple sets of 3D primitives [Liao et al. 2020], 3D meshes [Gao et al. 2022] and voxel grids [Nguyen-Phuoc et al. 2019] to radiance fields [Chan et al. 2020; Schwarz et al. 2021, 2022; Skorokhodov et al. 2022] and more recent Gaussian primitives [Barthel et al. 2024]. Some works generate a latent 3D representation that is rendered into a 2D feature map and then decoded into the final image [Gu et al. 2021; Niemeyer and Geiger 2021; Wewer et al. 2024]. This yields higher quality images, but at the cost of 3D inconsistencies across views.

Diffusion-based. Methods in this category are built upon denoising diffusion probabilistic models [Ho et al. 2020] to generate 3D assets. Akin to GAN-based models, works in this category include solutions that operate on different 3D representations, requiring direct 3D observations or indirect ones (e.g., 2D images). Methods that extract 3D information from 2D images in a pre-processing step before learning the diffusion model are sometimes referred to as two-stage approaches. Moreover, the diffusion model is either defined directly on the space of the target 3D representation, or inspired by Latent Diffusion Models [Rombach et al. 2022] on a latent space that is mapped to/from the 3D representation via learned decoder/encoder pairs. Among works that learn a 3D diffusion model from 3D observations (or have two stages), we find [Cai et al. 2020; Luo and Hu

2021] operating on 3D point clouds, [Müller et al. 2023] operating on grid-based radiance fields, and [Zhang et al. 2024] operating on 3D Gaussian primitives. Notably, [Chen et al. 2023] proposes a single-stage method that operates on the target 3D representation, namely tri-plane NeRF, but only requires indirect observations. Examples of methods operating on a latent space include [Li et al. 2023a; Zeng et al. 2022] for 3D point clouds and [Bautista et al. 2022; Ntavelis et al. 2023] for grid-based radiance fields. Similar to our work, [Li et al. 2023a] leverages a VQ-VAE [van den Oord et al. 2017] to construct the latent space, however, their focus is on the generation of object geometries as opposed to our latent 3D Gaussian diffusion enabling object and room-level view synthesis. There also exists a stream of works that use diffusion models directly in image space. Among those we have methods that optimize the 3D representation given the 2D supervision generated by a text-to-image diffusion model, like [Poole et al. 2022] for NeRFs and [Chen et al. 2024; Li et al. 2023b; Yi et al. 2024] for Gaussian primitives. However, these methods are limited to single object generation and their per-shape optimization approach is slower than our generations in a diffusion denoising process. We then find methods like [Anciukevičius et al. 2024] that denoise images by mapping them to the 3D representation and using rendering to map back to image space. In addition, there are works that use diffusion to generate 2D views of a hypothetical 3D scene directly [Liu et al. 2024; Watson et al. 2022]. These latter models can produce high-quality images, but they are potentially 3D inconsistent, and require typically some form of image conditioning, although unconditional generation could be achieved by pairing it with an unconditional image generator.

Our method falls into the category of diffusion-based models that operate in latent-space with Gaussian primitives as our underlying 3D representation. Following [Rombach et al. 2022], our model consists of a VQ-VAE that is trained on direct 3D observations to map to/from a latent representation and a diffusion model operating on the latter space. To our knowledge, we are the first method of this kind.

3 METHOD

We focus on the task of unconditional synthesis of 3D Gaussians primitives as a high-fidelity scene representation that features real-time rendering. To enable detailed 3D generation of objects and scalability to room-size scenes, our method lifts the 3D representation of Gaussian primitives to a learned, compressed latent space on which a diffusion model can efficiently operate. The generated latent representation is learned in a feature grid that can be decoded back to a set of 3D Gaussian primitives to support fast novel-view synthesis (Fig. 2). To efficiently map between the Gaussian primitives (Sec. 3.1) and the latent representation on which the diffusion model operates, we introduce a sparse convolutional network, which implements a VQ-VAE (Sec. 3.2). Finally, our latent diffusion model learns a denoising process in our low-dimensional latent space to unconditionally generate novel 3D Gaussian scenes from pure noise (Sec. 3.3).

3.1 Preliminaries: 3D Gaussian Splatting

Given a set of RGB images with camera poses, 3D Gaussian Splatting (3DG) [Kerbl et al. 2023] reconstructs the corresponding static scene, represented as a collection of 3D Gaussian primitives. Each Gaussian primitive comprises a 3D position $\mu_i \in \mathbb{R}^3$ and a 3D covariance matrix Σ_i that is factorized as $\Sigma_i := R_i S_i^2 R_i^\top$, where S_i is a nonnegative, diagonal scale matrix with diagonal denoted by $s_i \in \mathbb{R}_+$, and $R_i \in \text{SO}(3)$ is a rotation matrix represented as a unit quaternion $r_i \in \mathbb{R}^4$. To support rendering of RGB images, the view-dependent color $c_i(\gamma_i, d) \in \mathbb{R}^3$ of each Gaussian primitive is obtained from its spherical harmonics coefficients γ_i and the viewing direction d . In addition, each Gaussian primitive entails an opacity $\alpha_i \in \mathbb{R}_+$. An image C_π from camera π can be rendered by projecting and blending N depth-ordered Gaussians primitives as follows:

$$C_\pi(\mathbf{u}) := \sum_{i=1}^N c_i(\gamma_i, d) \omega_\pi^i(\mathbf{u}) \prod_{j=1}^{i-1} [1 - \omega_\pi^j(\mathbf{u})], \quad (1)$$

where $\omega_\pi^i(\mathbf{u}) := \alpha_i G_\pi^i(\mathbf{u})$ is the opacity of the i th primitive scaled by the contribution of the following function:

$$G_\pi^i(\mathbf{u}) := \exp \left[-\frac{1}{2} (\mathbf{u} - \mu_i^\pi)^\top (\Sigma_i^\pi)^{-1} (\mathbf{u} - \mu_i^\pi) \right]. \quad (2)$$

This represents the kernel of the 2D Gaussian with parameters $(\mu_i^\pi, \Sigma_i^\pi)$ that we obtain when projecting the primitive’s 3D Gaussian with parameters (μ_i, Σ_i) to the camera image plane under a linear approximation of the projection function (see [Kerbl et al. 2023] for more details).

The parameters of the 3D Gaussian primitives of a scene are optimized by minimizing an L_1 color loss \mathcal{L}_{RGB} and the negated structural similarity metric (SSIM) [Wang et al. 2004] between rendered images \hat{I} and target images I :

$$\mathcal{L}_{\text{3DG}} := (1 - \lambda_{\text{3DG}}) \mathcal{L}_{\text{RGB}} + \lambda_{\text{3DG}} (1 - \text{SSIM}(\hat{I}, I)), \quad (3)$$

where λ_{3DG} is a balancing factor.

3.2 Learning a Latent Space for 3D Gaussians

While 3D Gaussian Splatting offers an explicit representation that is highly expressive and efficient, the point cloud of 3D Gaussian primitives is spatially unstructured and sparse. The unstructured nature makes it challenging for a generalized model to learn. To recover spatial structure, we optimize primitives that are assigned to voxels of a sparse grid (Sec. 3.2.1). To cope with sparsity, we define a network comprising sparse convolutions to compress our 3D representation (Sec. 3.2.2) into a latent dense grid of low spatial resolution.

3.2.1 Sparse Grid-assigned 3D Gaussians. To train our VQ-VAE, we pre-compute for each scene 3D Gaussian primitives that are aligned with a sparse grid, i.e., the space of a scene is discretized into a 3D grid with voxel size d and each primitive is *uniquely* assigned to a voxel. This representation is optimized similar to [Kerbl et al. 2023] with a few differences. First, the position of a primitive μ_i is reparametrized in terms of a voxel index κ_i and a 3D displacement $\psi(\delta_i)$, depending on $\delta_i \in \mathbb{R}^3$, so that the primitive’s center becomes $\mu_i := \mathbf{y}_{\kappa_i} + \psi(\delta_i)$, where $\mathbf{y}_j \in \mathbb{R}^3$ denotes the 3D center of the j th voxel. Since each voxel can be assigned at most one Gaussian

primitive, the set of Gaussian primitives can be represented as a sparse grid θ of Gaussian parameters, where the κ_i th cell denoted by θ_{κ_i} contains the parameters of the i th Gaussian primitive, i.e. $\theta_{\kappa_i} := (\delta_i, \mathbf{s}_i, \mathbf{r}_i, \gamma_i, \alpha_i)$. During optimization, 3D primitives can move within their voxel and adjacent ones. To enforce this, we set $\psi(\delta) := 1.5 \tanh(\delta) d$. Second, we introduce a new densification strategy. A new Gaussian primitive is created in an inactive voxel, if an existing primitive from a neighboring cell moves into that voxel and the magnitude of its averaged view-space positional gradient exceeds a threshold ϵ_S , indicating the need for densification. The newly-created primitive is initialized at the center of the voxel (zero displacement) with isotropic scale d , identity rotation matrix, pre-defined small opacity and appearance averaged from all primitives competing for densification on the same voxel.

Akin to [Kerbl et al. 2023], primitives with opacity below a threshold ϵ_α are pruned and the proposed sparse, grid-assigned representation of 3D Gaussian primitives is optimized by minimizing \mathcal{L}_{3DG} as defined in Eq. (3).

3.2.2 3D Gaussian Compression Model. Our 3D Gaussian compression model is inspired by the success of latent diffusion [Rombach et al. 2022] for image synthesis. Specifically, we employ a VQ-VAE [van den Oord et al. 2017] due to its ability to learn an expressive prior over a small, discretized latent space, which is particularly valuable to handle the complexity of 3D space. Our network leverages 3D sparse convolutions to map between the sparse, grid-assigned Gaussians and a small latent dense grid. A vector quantization layer with a codebook of size K is employed at the bottleneck between encoder E and decoder D :

$$\mathbf{z}_e := E(\theta), \quad (4)$$

$$\mathbf{z}_q := \text{quantize}(\mathbf{z}_e), \quad (5)$$

$$\hat{\theta} := D(\mathbf{z}_q), \quad (6)$$

where \mathbf{z}_e is the output of the encoder given the sparse, grid-assigned 3D Gaussians θ , \mathbf{z}_q is the quantized sparse latent space and $\hat{\theta}$ is the reconstructed representation. \mathbf{z}_q serves as input to the diffusion model (Sec. 3.3), where it is converted to a low resolution dense grid. The 3D Gaussian compression network is trained with a VQ-VAE commitment loss $\mathcal{L}_{\text{commit}}$, to ensure the encoder commits to embeddings in the codebook. As reconstruction losses, we employ an L_1 color loss \mathcal{L}_{RGB} and a perceptual loss $\mathcal{L}_{\text{perc}}$ on M renderings of the reconstructed 3D Gaussians from different viewpoints. The perceptual loss encourages similarity of reconstructed and target features at different levels of detail.

The decoder employs generative sparse transpose convolutions [Gwak et al. 2020] to enable the generation of new coordinates in the upsampling. This is crucial to allow standalone usage of the decoder to decode synthesized latent grids from the diffusion model, without the possibility of leveraging cached coordinates from the encoder as in standard sparse transpose convolutions. The generation of new coordinates in each upsampling layer comes with the need for a pruning strategy to avoid an explosion in the number of active voxels. Thus, after each upsampling, a linear layer classifies each predicted voxel as occupied or free [Tatarchenko et al. 2017]. During training, these occupancies are supervised with a binary cross entropy loss (BCE) \mathcal{L}_{occ} using the grid-assigned 3D

Gaussians as target. At test time, they serve to effectively prune the predicted voxels. Thus, we define the combined loss function of the 3D Gaussian compression model as follows:

$$\mathcal{L}_{\text{comp}} := \lambda_{\text{commit}} \mathcal{L}_{\text{commit}} + \lambda_{\text{RGB}} \mathcal{L}_{\text{RGB}} + \lambda_{\text{perc}} \mathcal{L}_{\text{perc}} + \mathcal{L}_{\text{occ}}, \quad (7)$$

$$\text{where } \mathcal{L}_{\text{commit}} := \|\mathbf{z}_e - \mathbf{e}_\perp\|_2^2, \quad (8)$$

$$\mathcal{L}_{\text{perc}} := \|\Phi_{\text{VGG}}(\hat{I}) - \Phi_{\text{VGG}}(I)\|_2^2, \quad (9)$$

where \mathbf{e} are codebook entries and we use \perp to indicate that gradients to the embeddings are stopped. The codebook items are instead updated using an exponential moving average [van den Oord et al. 2017]. Φ_{VGG} is the vectorized concatenation of the first 5 feature layers before the max pooling operation of a VGG19 network [Simonyan and Zisserman 2015], where each layer is normalized by the square root of the number of elements.

The 3D Gaussian compression model learns a compact representation, where two downsampling layers of stride 2 lead to a volumetric compression by a factor of 64. At the same time, the number of parameters per voxel is drastically reduced to 4-element codebook items, where the codebook size is kept below 10k in all experiments.

3.3 Latent 3D Gaussian Diffusion

We propose a latent 3D diffusion model to learn the distribution of the compact latent space of 3D Gaussians p . To generate scenes from pure noise, without any prior knowledge, we need to use a dense grid in the latent diffusion, such that content may be synthesized anywhere in space. Hence, the compressed sparse grid is first converted to the corresponding low-resolution dense grid. To enable switching back to the sparse representation for decoding, the diffusion model is trained to denoise an additional occupancy element in the dense form.

The generation process is an inverse discrete-time Markov forward process. The forward process repeatedly adds Gaussian noise $\epsilon \in \mathcal{N}(0, I)$ to a sample $\mathbf{z}_0 \sim p$ leading to a series of increasingly noisy samples $\{\mathbf{z}_t | t \in [0, T]\}$. The noisy sample \mathbf{z}_t at a time step t is defined as

$$\mathbf{z}_t := \alpha_t \mathbf{z}_0 + \sigma_t \epsilon, \quad (10)$$

where parameters α_t and σ_t determine the amount of noise as part of the noise scheduling. After T noising steps the sample becomes pure Gaussian noise (i.e., $\alpha_T \approx 0$ and $\sigma_T \approx 1$). Our diffusion model reverses the forward process, i.e., it iteratively denoises a noisy sample beginning at T with pure noise, yielding at the end a clean sample \mathbf{z}_0 . We train our diffusion model $\hat{\mathbf{v}}_\phi(\mathbf{z}_t, t)$, parametrized by ϕ , to perform v-prediction [Salimans and Ho 2022], where the network output relates to the predicted clean sample $\hat{\mathbf{z}}_0$ by

$$\hat{\mathbf{z}}_0 := \alpha_t \mathbf{z}_t - \sigma_t \hat{\mathbf{v}}_\phi(\mathbf{z}_t, t). \quad (11)$$

This property is used to compute a mean squared error (MSE) loss against the noise-free sample to supervise the diffusion model:

$$\mathcal{L}_{\text{diff}} = \|\hat{\mathbf{z}}_0 - \mathbf{z}_0\|_2^2. \quad (12)$$

With our latent 3D Gaussian diffusion model the iterative synthesis process efficiently takes place in the low resolution latent space. A generated latent sample is sparsified using the occupancy channel,

and decoded to a high fidelity sparse 3D Gaussian representation using the decoder from Sec. 3.2.2.

3.4 Implementation Details

Sparse Grid-assigned 3D Gaussians. On all datasets, we scale the point cloud that is used to initialize the 3D Gaussians optimization into a unit cube and use a voxel size $d = 0.008$. While they use the same sparse grid resolution, scenes typically require $\sim 200k$ Gaussians whereas $\sim 8k$ are sufficient for objects. We set the densification and pruning thresholds $\epsilon_\delta = 0.0008$; $\epsilon_\alpha = 0.005$ and use $\lambda_{3DG} = 0.2$. View dependence is modeled with spherical harmonics of degree 1, which performs best across datasets. The optimization of 3D Gaussians requires $\sim 4\text{min}/\text{shape}$ and $\sim 10\text{min}/\text{room}$. Note that our implementation is unoptimized and significant improvement could be made through parallel processing.

3D Gaussian Compression Model. We use Minkowski Engine [Choy et al. 2019] to implement the 3D sparse convolutional network. Our convolutional blocks all use kernel size 3, with batch norm and ReLU activations. The encoder starts with a convolutional block increasing the input channels to 128. It is followed by two downsampling blocks, each consisting of two residual blocks, where the first doubles the number of channels, followed by a convolutional block with stride 2. Another residual block is employed in the bottleneck, where the number of channels is 512. A convolutional layer reduces the number of channels to 4 in the latent space, where the vector quantization is applied using codebook size $K = 4096$ on objects and $K = 8192$ on rooms. The decoder starts with a convolutional block that increases the number of channels to 512. This is followed by 2 upsampling blocks, each consisting of two residual blocks, where the first halves the number of channels, followed by a generative transpose convolution block [Gwak et al. 2020] with stride 2. After the upsampling, 2 residual blocks and a final convolutional layer map from 128 channels to the number of Gaussian parameters. We use loss weighting $\lambda_{\text{commit}} = 0.25$ on all datasets, $\lambda_{\text{RGB}} = 12.5$; $\lambda_{\text{perc}} = 0.1$ with $M = 4$ images on objects and $\lambda_{\text{RGB}} = 7.5$; $\lambda_{\text{perc}} = 0.3$ with $M = 12$ images on rooms. The batch size is 16 on objects and 4 on rooms. We train the model for 130/200/100 epochs on PhotoShape/ABO/3D-FRONT, using the Adam optimizer [Kingma and Ba 2015] with learning rates 0.0001/0.0002/0.0001 which are exponentially decayed by a factor of 0.998/0.98/0.95 at the end of each epoch. The training time on a single NVIDIA A100 GPU is $\sim 5\text{d}/1\text{d}/3.5\text{d}$ on PhotoShape/ABO/3D-FRONT.

3D Diffusion Model. The diffusion model is a 3D UNet, which adapts the architecture of [Dhariwal and Nichol 2021] to 3D. We use attention at resolutions 8 and 4 with 64 channels per head. The diffusion model operates on a 32^3 grid using a linear beta scheduling from 0.0001 to 0.02 in 1000 timesteps. We train the model for 1500/2500/3000 epochs on PhotoShape/ABO/3D-FRONT, using the Adam optimizer [Kingma and Ba 2015] with learning rate 0.0001 which is exponentially decayed by a factor of 0.998/0.9988/0.9988 at the end of each epoch. The training time on 2/1/1 NVIDIA A100 is $\sim 3.5\text{d}/1.5\text{d}/5\text{d}$ using batch size 64/32/16 on PhotoShape/ABO/3D-FRONT. For generation, we use DDPM sampling with 1000 steps.

4 EXPERIMENTS

In this section, we evaluate the performance of our method on unconditional generation of 3D assets. We also provide qualitative examples showcasing the ability of our method to generate room-scale scenes.

Datasets. We consider two benchmark datasets for the quantitative analysis, namely PhotoShape Chairs [Park et al. 2018] and Amazon Berkeley Objects (ABO) Tables [Collins et al. 2022]¹. Following [Müller et al. 2023], for PhotoShape Chairs, we consider 15,576 chairs rendered from 200 views along an Archimedean spiral. For ABO Tables, we use the provided 91 renderings from the upper hemisphere, considering 2-3 different environment map settings per object, resulting in 1676 tables split into 1520/156 for train/test. For room-scale scene generation, we train our model on ~ 2000 bedroom and living room style scenes from the 3D-FRONT [Fu et al. 2021] dataset. We render ~ 100 -500 images for training and ~ 20 -100 for testing, depending on the scene size. For all datasets, we use 512×512 images.

Metrics. We evaluate the quality of the generated 3D assets by measuring both the quality of the rendered images and their geometric plausibility. To evaluate the quality of the renderings, we use the Frechet Inception Distance [Heusel et al. 2018] (FID) and Kernel Inception Distance [Bińkowski et al. 2021] (KID) as implemented in [Obukhov et al. 2020]. All metrics are evaluated at 128×128 resolution. To evaluate the geometric plausibility, following [Achlioptas et al. 2018], we compute the Coverage Score (COV) and Minimum Matching Distance (MMD) using Chamfer Distance (CD), where the Coverage Score measures the diversity of the generated samples, while MMD assesses the quality of the generated samples. The geometry is extracted by voxelizing the 3D Gaussians and extracting a mesh using marching cubes, so that points on the surface can be sampled.

Baselines. We compare our method against state-of-the-art competitors that support unconditional generation of 3D assets and fall into both GAN-based and diffusion-based categories. Among GAN-based approaches, we consider π -GAN [Chan et al. 2020] and EG3D [Chan et al. 2022]. We also compare with the diffusion-based DiffRF [Müller et al. 2023]. All methods evaluated, including ours, use the same set of rendered images for training. GAN-based methods are trained directly on the rendered images, while DiffRF also uses per-shape radiance fields, which are pre-computed using the available posed training images. Similarly, our method uses pre-computed sparse grid-assigned 3D Gaussians (as per Section 3.2.1).

4.1 Baseline Comparison

As shown in Tabs. 1 and 2, our method leads to noticeable improvements compared to the baselines on all metrics. In particular, the perceptual metrics FID and KID show a large improvement, which indicates that our approach produces sharper, more detailed results. This is confirmed by the qualitative comparisons in Figs. 3 and 4. All compared approaches produce plausible shapes. However, by generating 3D Gaussian primitives, our method is able to synthesize

¹All objects for PhotoShape Chairs and ABO Tables were originally sourced from 3D Warehouse.



Fig. 3. Comparison on PhotoShape [Park et al. 2018]. Our method generates more detail than the baselines, such as thin structures, and has fewer artifacts.

Table 1. Quantitative comparison of unconditional generation on the PhotoShape Chairs [Park et al. 2018] dataset. MMD and KID scores are multiplied by 10^3 .

Method	FID ↓	KID ↓	COV ↑	MMD ↓
π -GAN [Chan et al. 2020]	52.71	13.64	39.92	7.387
EG3D [Chan et al. 2022]	16.54	8.412	47.55	5.619
DiffRF [Müller et al. 2023]	15.95	7.935	58.93	4.416
Ours	8.49	3.147	63.80	4.241

Table 2. Quantitative comparison of unconditional generation on the ABO Tables [Collins et al. 2022] dataset. MMD and KID scores are multiplied by 10^3 .

Method	FID ↓	KID ↓	COV ↑	MMD ↓
π -GAN [Chan et al. 2020]	41.67	13.81	44.23	10.92
EG3D [Chan et al. 2022]	31.18	11.67	48.15	9.327
DiffRF [Müller et al. 2023]	27.06	10.03	61.54	7.610
Ours	14.03	3.15	65.38	7.312
w/o compression model	197.1	166.8	51.92	8.483
w/o RGB loss	17.26	3.28	63.46	7.756
w/o perceptual loss	34.35	13.65	61.53	7.488

thinner structures, such as chair and table legs, where the DiffRF results are more coarse due to the limiting radiance field grid resolution. The GAN-based approach EG3D shows more artifacts and view-dependent inconsistencies, e.g., in the chair leg areas.

Tab. 3 provides a runtime comparison. By synthesizing 3D Gaussians, which can be very efficiently rasterized, our method achieves significantly faster rendering speed, i.e., ~ 50 times faster than DiffRF using radiance fields. The GAN-based EG3D generates shapes in a single network forward pass, hence has much faster generation time compared to diffusion-based approaches. Nonetheless, our method almost halves generation time compared to DiffRF.

4.2 Ablation Study

To verify design choices of our method, we perform an ablation study on the ABO Tables dataset. The quantitative evaluation in Tab. 2, as well as the qualitative comparison in Fig. 5 demonstrate that the full version of our method leads to the best performance.

Without compression model. Omitting the 3D Gaussian compression model, i.e., training the diffusion model directly on optimized, grid-assigned 3D Gaussians (Sec. 3.2.1) of low resolution (32^3), results in a drastic performance drop. We found that the diffusion model struggles to denoise this more complex, higher dimensional



Fig. 4. Comparison on ABO [Collins et al. 2022]. Tables generated by our method are sharper and show less artifacts compared to the baselines.



Fig. 5. Ablation study on ABO [Collins et al. 2022]. Training the 3D Gaussian compression model without the rendering losses (perceptual or RGB) leads to more blurry results, especially without perceptual loss. The variant without RGB loss additionally produces less color variations in the generated scenes.



Fig. 6. Qualitative results on unconditional room generation on 3D-FRONT [Fu et al. 2021]. Our method scales to room-size scenes and synthesizes plausible geometry and appearance. We visualize the generated 3D Gaussian ellipsoids and their renderings.

Table 3. Runtime comparison on ABO Tables [Collins et al. 2022] dataset using one NVIDIA RTX A6000. By generating 3D Gaussians, our method enables much faster rendering speed. With a single forward pass, EG3D generates faster than diffusion-based approaches.

Method	Generation time ↓ per shape	Rendering time ↓ per frame	
EG3D [Chan et al. 2022]	6ms	23ms	@ 128 × 128
DiffRF [Müller et al. 2023]	21s	48ms	@ 512 × 512
Ours	13s	0.91ms	@ 512 × 512

space of 3D Gaussian parameters θ_{K_i} , compared to our vector-quantized latent features of 4 elements, where the codebook size allows for less than 10k alternative embeddings. After the same training time as our full method, the version “w/o compression model” still struggles to generate Gaussians that coherently describe the 3D shape (see examples in Fig. 7).



Fig. 7. Ablation experiment “w/o compression model” struggles to generate coherent 3D Gaussians after the same training time as our complete method.

Without RGB. Dropping the RGB loss during training of the 3D Gaussian compression model leads to reduced perceptual and geometric metrics. Qualitatively, we observe that without the RGB loss the renderings tend to be more blurry and the color variety of the generated shapes seems reduced (Fig. 5).

Without perceptual loss. The experiment without perceptual loss in the 3D Gaussian compression training shows a clear decrease in the performance measured by all metrics. The renderings lose sharpness, e.g., the wooden patterns on the tables in Fig. 5 are no longer visible.

4.3 Other Qualitative Results

4.3.1 Unconditional Scene Generation. We showcase the ability of our latent 3D Gaussian diffusion to scale to room-size scenes. The sparse 3D Gaussians compression model enables flexible scaling to operate on scenes, which have ~200k Gaussians, compared to ~8k on objects, while the diffusion model can still operate on the same latent space dimension as for the object-level datasets. Fig. 6 shows results on unconditional generation of rooms, where the model is trained on bedroom and living room style scenes from the 3D-FRONT [Fu et al. 2021] dataset. The generated scenes have plausible and varied configurations of furniture, and an accurate geometry, which is visible in the visualization of 3D Gaussian ellipsoids.

4.3.2 Nearest Neighbors in the Training Set. Fig. 8 visualizes our generated chairs next to their nearest neighbors from the training set of optimized sparse grid-assigned 3D Gaussians. The geometric nearest neighbors are computed using Chamfer Distance on point

clouds sampled from the 3D Gaussians. We observe that the generated chairs are substantially different from their nearest neighbors, indicating that the model does not purely retrieve from the training set, but generates novel shapes.



Fig. 8. Visualization of geometric nearest neighbors in the training set using Chamfer Distance. Our approach can generate novel samples (left) that are different from their nearest neighbors in the training set (right).

4.4 Limitations

While our method is among the first to show its applicability to 3D scene generation at room-scale, we believe there are still significant open challenges. One key ingredient towards achieving the scalability of our approach lies in the latent 3D scene representation of the 3D Gaussians. Here, analog to 2D image diffusion models [Rombach et al. 2022], larger neural network models will facilitate the creation of outputs of larger scene extents and higher visual fidelity. In this context, available computational resources was a major bottleneck that limited further exploration. However, at the same time, we believe that additional training strategies, e.g., exploiting spatial subdivision strategies of 3D spaces, could further alleviate memory and computational limitations.

At the same time, our method is currently trained on synthetic datasets such as PhotoShape [Park et al. 2018], ABO [Collins et al. 2022], or 3D-FRONT [Fu et al. 2021]. Here, we can see a future potential on real-world datasets that provide ground truth 3D supervision at the scene level. Unfortunately, 3D datasets with high-fidelity DSLR captures (which is required to reconstruct the Gaussian ground truth pairs), such as Tanks and Temples [Knapitsch

et al. 2017] or ScanNet++ [Yeshwanth et al. 2023] are still relatively limited in terms of the number of available 3D scenes.

5 CONCLUSION

We have presented L3DG, a novel generative approach that models a 3D scene distribution represented by 3D Gaussians. The core idea of our method is a latent 3D diffusion model whose latent space is learned by a VQ-VAE for which we propose a sparse convolutional 3D architecture. This facilitates the scalability of our method and significantly improves the visual quality over existing works. For instance, in comparison to NeRF-based generators, such as DiffRF [Müller et al. 2023], L3DG can be rendered faster and thus trained on larger scenes. In particular, this allows us to showcase a first step towards room-scale scene generation. Overall, we believe that our method is an important stepping stone to support the 3D content generation process along a wide range of applications in computer graphics.

ACKNOWLEDGMENTS

This work was funded by a Meta SRA. Matthias Nießner was also supported by the ERC Starting Grant Scan2CAD (804724) and Angela Dai was supported by the ERC Starting Grant SpatialSem (101076253).

REFERENCES

- Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. 2018. Learning Representations and Generative Models for 3D Point Clouds. arXiv:1707.02392 [cs.CV]
- Titas Aciukevičius, Zexiang Xu, Matthew Fisher, Paul Henderson, Hakan Bilen, Niloy J. Mitra, and Paul Guerrero. 2024. RenderDiffusion: Image Diffusion for 3D Reconstruction, Inpainting and Generation. arXiv:2211.09869 [cs.CV]
- Florian Barthel, Arian Beckmann, Wieland Morgenstern, Anna Hilsmann, and Peter Eisert. 2024. Gaussian Splattting Decoder for 3D-aware Generative Adversarial Networks. arXiv:2404.10625 [cs.CV]
- Miguel Angel Bautista, Pengsheng Guo, Samira Abnar, Walter Talbott, Alexander Toshev, Zhuoyuan Chen, Laurent Dinh, Shuangfei Zhai, Hanlin Goh, Daniel Ulbricht, Afshin Dehghan, and Josh Susskind. 2022. GAUDI: A Neural Architect for Immersive 3D Scene Generation. arXiv:2207.13751 [cs.CV]
- Mikolaj Binkowski, Danica J. Sutherland, Michael Arbel, and Arthur Gretton. 2021. Demystifying MMD GANs. arXiv:1801.01401 [stat.ML]
- Ruojin Cai, Guandao Yang, Hadar Averbuch-Elor, Zekun Hao, Serge Belongie, Noah Snavely, and Bharath Hariharan. 2020. Learning Gradient Fields for Shape Generation. arXiv:2008.06520 [cs.CV]
- Eric Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. 2020. pi-GAN: Periodic Implicit Generative Adversarial Networks for 3D-Aware Image Synthesis. In *arXiv*.
- Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. 2022. Efficient Geometry-aware 3D Generative Adversarial Networks. In *CVPR*.
- Hansheng Chen, Jiatao Gu, Anpei Chen, Wei Tian, Zhiwen Tu, Lingjie Liu, and Hao Su. 2023. Single-Stage Diffusion NeRF: A Unified Approach to 3D Generation and Reconstruction. In *ICCV*.
- Zilong Chen, Feng Wang, Yikai Wang, and Huaping Liu. 2024. Text-to-3D using Gaussian Splattting. arXiv:2309.16585 [cs.CV]
- Christopher Choy, JunYoung Gwak, and Silvio Savarese. 2019. 4D Spatio-Temporal ConvNets: Minkowski Convolutional Neural Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3075–3084.
- Jasmine Collins, Shubham Goel, Kenan Deng, Achleshwar Luthra, Leon Xu, Erhan Gundogdu, Xi Zhang, Tomas F Yago Vicente, Thomas Dideriksen, Himanshu Arora, et al. 2022. Abo: Dataset and benchmarks for real-world 3d object understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 21126–21136.
- Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems* 34 (2021). 8780–8794.
- Huan Fu, Bowen Cai, Lin Gao, Ling-Xiao Zhang, Jiaming Wang, Cao Li, Qixun Zeng, Chengyue Sun, Rongfei Jia, Binqiang Zhao, et al. 2021. 3d-front: 3d furnished rooms with layouts and semantics. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 10933–10942.
- Jun Gao, Tianchang Shen, Zian Wang, Wenzheng Chen, Kangxue Yin, Daiqing Li, Or Litany, Zan Gojcic, and Sanja Fidler. 2022. GET3D: A Generative Model of High Quality 3D Textured Shapes Learned from Images. arXiv:2209.11163 [cs.CV]
- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative Adversarial Networks. arXiv:1406.2661 [stat.ML]
- Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. 2021. StyleNeRF: A Style-based 3D-Aware Generator for High-resolution Image Synthesis. arXiv:2110.08985 [cs.CV]
- JunYoung Gwak, Christopher B Choy, and Silvio Savarese. 2020. Generative Sparse Detection Networks for 3D Single-shot Object Detection. In *European conference on computer vision*.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2018. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. arXiv:1706.08500 [cs.LG]
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising Diffusion Probabilistic Models. arXiv:2006.11239 [cs.LG]
- James T Kajiya and Brian P Von Herzen. 1984. Ray tracing volume densities. *ACM SIGGRAPH computer graphics* 18, 3 (1984), 165–174.
- Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023. 3D Gaussian Splattting for Real-Time Radiance Field Rendering. *ACM Transactions on Graphics* 42, 4 (July 2023). <https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/>
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. *CoRR* abs/1412.6980 (2015).
- Arno Knapsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. 2017. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)* 36, 4 (2017), 1–13.
- Xinhai Li, Huaibin Wang, and Kuo-Kun Tseng. 2023b. GaussianDiffusion: 3D Gaussian Splattting for Denoising Diffusion Probabilistic Models with Structured Noise. arXiv:2311.11221 [cs.CV]
- Yuhan Li, Yishun Dou, Xuanhong Chen, Bingbing Ni, Yilin Sun, Yutian Liu, and Fuzhen Wang. 2023a. 3DQD: Generalized Deep 3D Shape Prior via Part-Discretized Diffusion Process. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Yiyi Liao, Katja Schwarz, Lars Mescheder, and Andreas Geiger. 2020. Towards Unsupervised Learning of Generative Models for 3D Controllable Image Synthesis. arXiv:1912.05237 [cs.CV]
- Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. 2024. SyncDreamer: Generating Multiview-consistent Images from a Single-view Image. arXiv:2309.03453 [cs.CV]
- Shitong Luo and Wei Hu. 2021. Diffusion Probabilistic Models for 3D Point Cloud Generation. arXiv:2103.01458 [cs.CV]
- Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2020. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*. Springer, 405–421.
- Norman Müller, Yawar Siddiqui, Lorenzo Porzi, Samuel Rota Bulò, Peter Kortscheder, and Matthias Nießner. 2023. Diffr: Rendering-guided 3d radiance field diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4328–4338.
- Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. 2019. HoloGAN: Unsupervised learning of 3D representations from natural images. arXiv:1904.01326 [cs.CV]
- Michael Niemeyer and Andreas Geiger. 2021. GIRAFFE: Representing Scenes as Compositional Generative Neural Feature Fields. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*.
- Evangelos NTavelis, Aliaksandr Siarohin, Kyle Olszewski, Chaoyang Wang, Luc Van Gool, and Sergey Tulyakov. 2023. AutoDecoding Latent 3D Diffusion Models. arXiv:2307.05445 [cs.CV]
- Anton Obukhov, Maximilian Seitzer, Po-Wei Wu, Semen Zhydenko, Jonathan Kyl, and Elvis Yu-Jing Lin. 2020. High-fidelity performance metrics for generative models in PyTorch. <https://doi.org/10.5281/zenodo.4957738> Version: 0.3.0, DOI: 10.5281/zenodo.4957738.
- Keunhong Park, Konstantinos Rematas, Ali Farhadi, and Steven M. Seitz. 2018. Photo2Shape: Photorealistic Materials for Large-Scale Shape Collections. *ACM Trans. Graph.* 37, 6, Article 192 (Nov. 2018).
- Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. 2022. DreamFusion: Text-to-3D using 2D Diffusion. *arXiv* (2022).
- R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. 2022. High-Resolution Image Synthesis with Latent Diffusion Models. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE Computer Society, Los Alamitos, CA, USA, 10674–10685. <https://doi.org/10.1109/CVPR52688.2022.01042>
- Tim Salimans and Jonathan Ho. 2022. Progressive Distillation for Fast Sampling of Diffusion Models. In *International Conference on Learning Representations*.

- Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. 2021. GRAF: Generative Radiance Fields for 3D-Aware Image Synthesis. arXiv:2007.02442 [cs.CV]
- Katja Schwarz, Axel Sauer, Michael Niemeyer, Yiyi Liao, and Andreas Geiger. 2022. VoxGRAF: Fast 3D-Aware Image Synthesis with Sparse Voxel Grids. arXiv:2206.07695 [cs.CV]
- Karen Simonyan and Andrew Zisserman. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1409.1556>
- Ivan Skorokhodov, Sergey Tulyakov, Yiqun Wang, and Peter Wonka. 2022. EpiGRAF: Rethinking training of 3D GANs. arXiv:2206.10535 [cs.CV]
- M. Tatarchenko, A. Dosovitskiy, and T. Brox. 2017. Octree Generating Networks: Efficient Convolutional Architectures for High-resolution 3D Outputs. In *IEEE International Conference on Computer Vision (ICCV)*. <http://lmb.informatik.uni-freiburg.de/Publications/2017/TDB17b>
- Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. 2017. Neural discrete representation learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems* (Long Beach, California, USA) (NIPS’17). Curran Associates Inc., Red Hook, NY, USA, 6309–6318.
- Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* 13, 4 (2004), 600–612.
- Daniel Watson, William Chan, Ricardo Martin-Brualla, Jonathan Ho, Andrea Tagliasacchi, and Mohammad Norouzi. 2022. Novel View Synthesis with Diffusion Models. arXiv:2210.04628 [cs.CV]
- Christopher Wewer, Kevin Raj, Eddy Ilg, Bernt Schiele, and Jan Eric Lenssen. 2024. latentSplat: Autoencoding Variational Gaussians for Fast Generalizable 3D Reconstruction. arXiv:2403.16292 [cs.CV]
- Chandan Yeshwanth, Yueh-Cheng Liu, Matthias Nießner, and Angela Dai. 2023. ScanNet++: A High-Fidelity Dataset of 3D Indoor Scenes. In *Proceedings of the International Conference on Computer Vision (ICCV)*.
- Taoran Yi, Jiemin Fang, Junjie Wang, Guanjun Wu, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Qi Tian, and Xinggang Wang. 2024. GaussianDreamer: Fast Generation from Text to 3D Gaussians by Bridging 2D and 3D Diffusion Models. arXiv:2310.08529 [cs.CV]
- Xiaohui Zeng, Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany, Sanja Fidler, and Karsten Kreis. 2022. LION: Latent Point Diffusion Models for 3D Shape Generation. arXiv:2210.06978 [cs.CV]
- Bowen Zhang, Yiji Cheng, Jiaolong Yang, Chunyu Wang, Feng Zhao, Yansong Tang, Dong Chen, and Baining Guo. 2024. GaussianCube: Structuring Gaussian Splatting using Optimal Transport for 3D Generative Modeling. arXiv:2403.19655 [cs.CV]