

TexGaussian: Generating High-quality PBR Material via Octree-based 3D Gaussian Splatting

Bojun Xiong^{1*†}, Jialun Liu^{2*}, Jiakui Hu^{3†}, Chenming Wu², Jinbo Wu², Xing Liu²,
Chen Zhao², Errui Ding², Zhouhui Lian^{1‡}

¹Wangxuan Institute of Computer Technology, Peking University

²Baidu VIS

³Institute of Medical Technology, Peking University

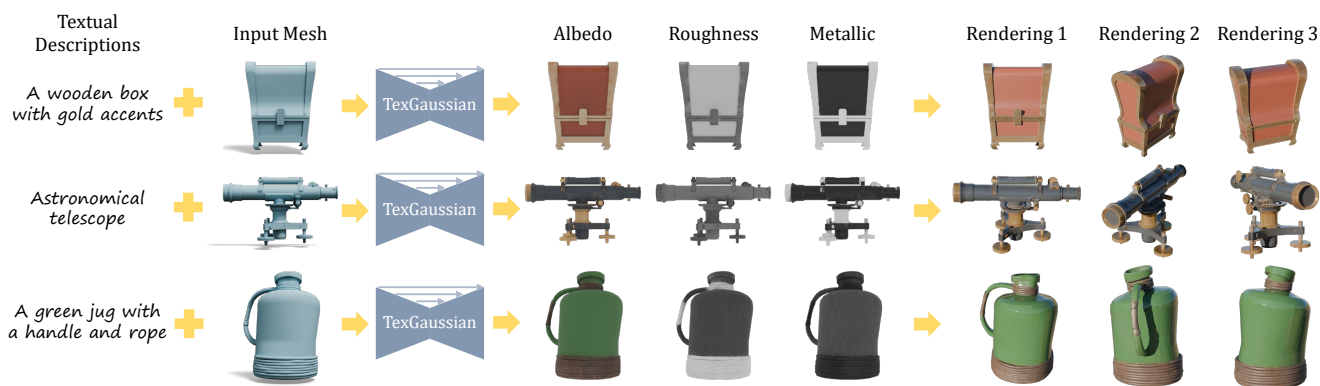


Figure 1. Our proposed TexGaussian is capable of generating high-quality PBR material given the input 3D mesh based on the corresponding textual descriptions. The generated results are naturally compatible with modern graphical engines for photo-realistic rendering under different environment maps.

Abstract

Physically Based Rendering (PBR) materials play a crucial role in modern graphics, enabling photorealistic rendering across diverse environment maps. Developing an effective and efficient algorithm that is capable of automatically generating high-quality PBR materials rather than RGB texture for 3D meshes can significantly streamline the 3D content creation. Most existing methods leverage pre-trained 2D diffusion models for multi-view image synthesis, which often leads to severe inconsistency between the generated textures and input 3D meshes. This paper presents *TexGaussian*, a novel method that uses octant-aligned 3D Gaussian Splatting for rapid PBR material generation. Specifically, we place each 3D Gaussian on the finest leaf node of the octree built from the input 3D mesh

to render the multi-view images not only for the albedo map but also for roughness and metallic. Moreover, our model is trained in a regression manner instead of diffusion denoising, capable of generating the PBR material for a 3D mesh in a single feed-forward process. Extensive experiments on publicly available benchmarks demonstrate that our method synthesizes more visually pleasing PBR materials and runs faster than previous methods in both unconditional and text-conditional scenarios, which exhibit better consistency with the given geometry. Our code and trained models are available at <https://3d-aigc.github.io/TexGaussian>.

1. Introduction

Traditional 3D asset creation relies heavily on the expertise and extensive effort of professional designers [22], posing a significant barrier for casual users interested in creating 3D models independently. In the 3D design process, geometry creation typically represents only a small por-

*Denotes equal contribution.

†This work was partly done when Bojun Xiong and Jiakui Hu interned in Baidu VIS.

‡Corresponding author. E-mail: lianzhouhui@pku.edu.cn

tion of the overall time, while the majority is dedicated to developing textures and appearances, which are far more time-consuming. Achieving a delicate appearance for a 3D model often demands substantial time and effort from experienced designers.

Recently, Artificial Intelligence Generated Content (AIGC) based on deep generative models, especially diffusion models [16, 45] have been widely used to facilitate the process of artistic creation, catalyzing advancements in image generation [14, 36, 39, 41, 42] and video generation [2–4, 17, 51]. As a result, exploring effective ways to leverage deep generative models to streamline the creation of detailed appearances for 3D models has become a popular direction in the graphics and vision communities.

Recent advancements in 3D texture generation attempt to use depth-conditional pre-trained 2D diffusion model [56] to synthesize RGB images based on the depth maps rendered from multiple views, such as TEXTure [40] and text2tex [7]. Subsequent works [5, 19, 55] further improve multi-view consistency via blending the multi-view images to a single and consistent texture map in every denoising step. However, these methods struggle to have a global picture of 3D geometries due to the use of 2D diffusion models, leading to inconsistencies between the texture map and the semantics of the input 3D meshes. Moreover, the generated assets suffer from illumination-baked textures, which can significantly degrade the quality of the final rendering when placed in novel lighting conditions [54]. While DreamMat [58] supplements geometry and light control to achieve material decomposition through score distillation sampling [37], it still struggles to fully capture the global geometry features. This limitation often results in the multi-face Janus problem and leads to over-saturated colors.

On the other hand, training a 3D neural network directly on 3D data, such as Point-UV Diffusion [53] and TexOct [25], is an effective way for 3D global consistency. Meanwhile, this avoids multi-view sampling and score distillation sampling, which accelerates the process of texture synthesis. However, relying on colored point clouds for 3D representation and supervision often results in blurred outputs, primarily due to the sparse and non-compact nature of point clouds in 3D space. Due to the limitations of the adopted 3D representations and the lack of training data on PBR materials, these approaches are incapable of generating high-fidelity PBR materials for 3D models.

To address the aforementioned challenges, this paper presents TexGaussian, a fast and high-fidelity PBR material generation model directly in 3D space that maintains 3D global consistency. Different from previous approaches that primarily rely on diffusion models, our method works in a regression manner to regress the PBR material from the input mesh for faster generation speed. To enable effective learning in 3D space, we propose to use octree, a

specialized sparse voxel structure that efficiently organizes and preserves 3D information, which can be built from 3D point clouds sampled from the surface of the object. However, directly regressing the color of 3D point clouds on octree often results in blurry textures as mentioned in [25]. To tackle the challenges of the incompact and discrete nature presented by points, we use 3D Gaussian Splatting (3DGS) [21], a robust representation that bridges the gap between 3D space and 2D raster images, allowing us to fully utilize rich 2D image information to alleviate blurring results. Specifically, for each input mesh, we sample the dense 3D point clouds on its surface to build the corresponding octree. On each octant (i.e., the finest leaf node of octree), we place a 3D Gaussian [21] at its central position. Then, we use the octree-based 3D U-Net [48] to predict the parameters of each 3D Gaussian on octants. Apart from RGB colors, we extend each 3D Gaussian with additional parameters to represent the roughness and metallicity of 3D objects. Multi-view images, including albedo, roughness, and metallic maps, can be rasterized from all these 3D Gaussians via 3DGS. The 3D U-Net is supervised by the difference between the predicted multi-view images and their corresponding ground truth. Notably, the 3D U-Net is trained to directly regress the multi-view images based on the geometry feature of the input 3D model, which further facilitates the process of 3D Gaussian prediction compared to diffusion manner. We train our TexGaussian model on a subset of Objaverse [12] with high-quality PBR materials, enabling fast PBR material generation with a single feed-forward pass. In summary, the contributions of our paper are threefold:

- We propose an octant-aligned 3D Gaussian Splatting method for high-quality PBR material synthesis on untextured input 3D mesh, which fully utilizes the supervision from 2D images, avoiding blurry results caused by the discreteness of 3D point clouds.
- We adopt a regression manner to train our 3D U-Net model instead of diffusion denoising, achieving faster generative speed.
- We propose TexGaussian, a novel PBR material generation method based on the above two techniques. To our knowledge, our method first generates PBR material directly in 3D space. Qualitative and quantitative experiments have been conducted to verify the superiority of the quality and efficiency of our method over other existing approaches.

2. Related Work

In this section, we mainly summarize current texture synthesis methods, which can be roughly divided into three categories.

2.1. Multi-view Images Synthesis

Many previous works have tried to leverage the powerful T2I model to assist texture generation for 3D shapes. Specifically, they render the depth map of input 3D mesh from multiple views and use depth conditional T2I models [56] to synthesize RGB images and perform text-conditioned texture synthesis. TEXTure [40] and Text2Tex [7] iteratively paint a mesh from different views. However, images synthesized in early view could produce errors that are not reconcilable with the geometry that is observed in later views. Many subsequent works try to alleviate multi-view inconsistency via different alignment modules. TexFusion [5] proposes a sequential interlaced multi-view sampler that interleaves texture assembling with denoising steps in different camera views. Similarly, TexGen [19] directly enforces view consistent sampling in RGB texture space and develops a noise resampling strategy to retain rich texture details. TexPainter [55] blends images from different views into a common color-space texture image by weighted averaging to guarantee multi-view consistency. GenesisTex [15] introducing style consistency and dynamic alignment across multiple viewpoints. To remove light influence from 2D diffusion models, Paint3D [54] contribute separate UV Inpainting and UVHD diffusion models specialized in shape-aware refinement. Although these methods achieve impressive texture results, they can still hardly comprehend the overall geometry of input 3D mesh.

2.2. Optimization-based 3D Generation

Before the emergence of large-scale Text-to-Image generative model, earlier methods [9, 18, 29, 30, 32] propose to optimize texture map of 3D object via natural language supervised visual model, CLIP [38]. Subsequently, score distillation sampling (SDS) was adopted by DreamFusion [37] and Magic3D [24]. The key idea is to optimize 3D representations such as NeRF [31] or InstantNGP [33] with the gradient guidance from 2D diffusion priors [37, 50]. To generate PBR material, TextureDreamer [52] optimizes spatially-varying bidirectional reflectance distribution (BRDF) field through personalized geometric-aware score distillation. Fantasia3D [8] uses a single predefined environmental. However, the generated images from diffusion models may not be consistent with the given environment light. DreamMat [58] proposes a novel geometry and light-aware diffusion model, which is trained to generate images that are consistent with the given environment light. FlashTex[13] also proposes a light-conditioned diffusion model within a two-stage pipeline, combining reconstruction and SDS optimization to enhance texture quality and achieve better light disentanglement. However, these methods struggle with the Janus problem due to the semantically ambiguous. And the time consumption is relatively too long to use in practice.

2.3. Generating Texture from 3D Data

The most straightforward way to synthesize texture map for 3D mesh is to train generative model directly from 3D data with texture ground truth [6, 10–12]. Early methods such as Texture Fields [34] learn implicit texture fields to assign a color to each pixel on the surface of the 3D shape. Texturify [44] devices face convolution operation on mesh surface to predict texture on each face. It employs differentiable rendering with an adversarial loss to ensure that generated textures produce realistic imagery. Recently, some diffusion-based texture synthesis methods, such as PointUV [53] and TexOct [25] train a denoising network on colors of point clouds which are further mapped to 2D UV map. Although these methods achieve better 3D global consistency with input mesh, they are only trained on several categories of small datasets [6]. What’s more, discrete supervision from 3D point clouds leads to suboptimal results compared with continuous signals such as 2D images.

3. Method

In this section, we provide a detailed explanation of our proposed method, TexGaussian. The overall pipeline of our method is shown in Fig 2. Existing texture synthesis approaches mainly rely on pre-trained 2D diffusion models, which struggle to fully understand the overall 3D structure. This often leads to misalignment between the generated texture map and 3D semantics. Our goal is to synthesize high-quality PBR materials for a given mesh directly in 3D space.

3.1. Overview

To enable effective learning in 3D space, we use octree, a sparse voxel structure, to organize and store 3D information without compromising representation quality due to the inherent sparseness of 3D objects in 3D space. Thus, we sample a large number of points on the surface of the given mesh to construct the corresponding octree. The key components of our method are the octant-aligned Gaussian Splatting and the octree-based 3D U-Net. Specifically, we place a 3D Gaussian at the center of each octree’s finest leaf node. The octree-based 3D-U-net is trained to predict the parameters of each 3D Gaussian. Under this circumstance, the generated octant-aligned 3D Gaussians are naturally on the mesh surface. We render them from multiple viewpoints using 3D Gaussian Splatting and train the 3D U-Net by minimizing the difference in 2D raster images and 3D Gaussian parameters. During inference, the rendered multi-view images are baked into the UV space of input mesh using differentiable mesh rendering, producing the final texture and material map. The details of 3D Gaussian Splatting (3DGS) are provided in the supplementary material.

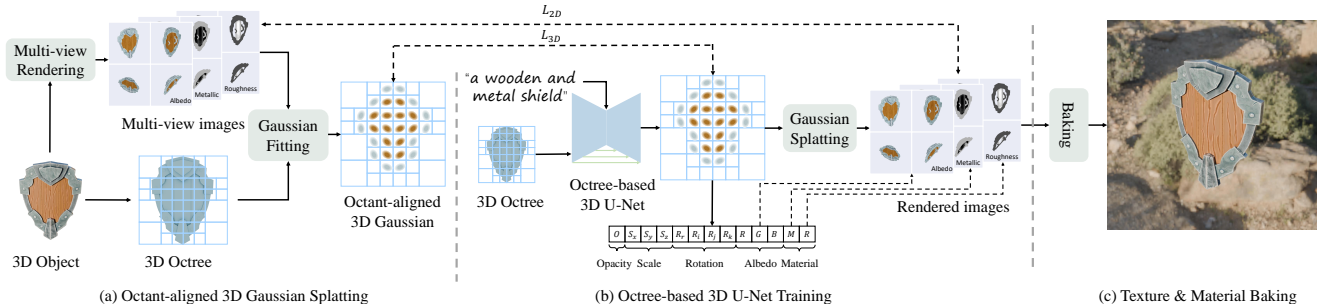


Figure 2. An overview of our PBR material generation framework. (a) We propose octant-aligned 3D Gaussian Splatting, which positions a 3D Gaussian at the center of each finest leaf node of the constructed octree. Additional channels are added at the end of the Gaussian parameters to model PBR material. (b) We use the 3D U-Net built upon octree-based convolutional networks to predict the Gaussian parameters. Our octree-based 3D U-Net is trained by minimizing the difference on 2D raster images and 3D Gaussian parameters. (c) We bake the multi-view rendered images to the UV space of the input 3D model to realize physically based rendering under new illumination environments.

3.2. Octant-aligned 3D Gaussian Splatting

For a given mesh, we first sample $N = 100,000$ 3D points on its surface. The corresponding octree is then built by adaptively subdividing the voxels containing those points until the maximum depth is reached. As a result, all of the finest leaf nodes of octree lie along the boundary of the 3D object, which is consistent with the characteristic of the optimized 3D Gaussian. Therefore, we align a 3D Gaussian at the central position of each finest leaf node to effectively model the appearance without compromising the splatting quality. It is worth noting that we do not adjust the position of 3D Gaussians because they are already on the surface of 3D shape and adding additional offset relative to the octant center would not improve the rendering quality through our early experiments. To model PBR material that includes roughness and metallic information, we follow [20, 43] to append two additional channels at the end of 3D Gaussian parameters which are responsible for roughness and metallic map rendering, respectively.

For each 3D object in our dataset, we render multi-view images of albedo, roughness, and metallic maps for the training of 3D Gaussian. The multi-view rendering results of the albedo map are view-independent. So we just use three RGB channels to take the place of the original spherical harmonics. As noted, we exclude the position from 3D Gaussian parameters in our model. Thus, our 3D Gaussian parameters consist of 13 channels in total: one for opacity, three for scale, four for rotation, three for albedo, one for roughness, and one for metallic. To stabilize training, we choose to employ different activation functions compared to the original Gaussian Splatting [21]. Specifically, We multiply the softplus-activated scales s_i with 0.01, ensuring that the initial 3D Gaussians conform to the object’s counter of object at the beginning of training rather than expanding outward.

We pre-fitting the parameters of 3D Gaussians on the

constructed octree for each 3D object in our dataset via the original loss function in [21] on multi-view RGB images, roughness maps, and metallic maps:

$$L = (1 - \lambda)L_1 + \lambda L_{D-SSIM}, \quad (1)$$

where $\lambda = 0.2$.

3.3. Octree-based 3D U-Net Training

To handle the encoding of our octant-aligned 3D Gaussian representation, we use the 3D U-Net built upon octree-based convolutional neural networks [48] to predict the Gaussian parameters. Inspired by LGM [46], the output feature of the 3D U-Net on each octant is treated as the 3D Gaussian parameters, which contain 13 channels, as discussed in the last subsection.

To effectively train our 3D U-Net, we adopt the regressive loss objective, which could further facilitate the generation process. The input to our octree-based 3D U-Net is the geometry feature on each octree’s finest leaf node, such as normal and local displacement. For text-conditioned PBR material synthesis, the text feature is extracted by pre-trained CLIP model [38] and is fed to U-Net via the octree-based multi-head cross attention mechanism similar to [47]. The predicted 3D Gaussians are rasterized from multiple views via Gaussian Splatting [21]. At each training step, we rasterize the RGB images, alpha images, roughness and metallic maps from randomly selected eight views. Following [46], we apply mean square error (MSE) loss and VGG-based LPIPS loss [57] to the RGB image, roughness map, and metallic map:

$$L_{RGB} = L_{MSE}(I_{RGB}, I_{RGB}^{GT}) + L_{LPIPS}(I_{RGB}, I_{RGB}^{GT}), \quad (2)$$

$$L_R = L_{MSE}(I_R, I_R^{GT}) + L_{LPIPS}(I_R, I_R^{GT}), \quad (3)$$

$$L_M = L_{MSE}(I_M, I_M^{GT}) + L_{LPIPS}(I_M, I_M^{GT}), \quad (4)$$

where ‘R’ and ‘M’ denote roughness and metallic, respectively. We further apply the MSE loss on the alpha image for faster convergence of the shape:

$$L_{\alpha} = L_{\text{MSE}}(I_{\alpha}, I_{\alpha}^{\text{GT}}). \quad (5)$$

To accelerate the coverage process, we also apply the 3D MSE loss $L_{3\text{D}}$, which calculates the difference between predicted parameters of 3D Gaussians and pre-fitting ones. Finally, the complete loss function of our model is defined as the sum of all the above losses:

$$L_{\text{total}} = L_{\text{RGB}} + L_{\text{R}} + L_{\text{M}} + L_{\alpha} + L_{3\text{D}}. \quad (6)$$

3.4. Texture and Material Baking

In the inference stage, we also first build the corresponding octree for the input 3D mesh. Then, we use our trained octree-based 3D U-Net to generate 3D Gaussian on every octant and rasterize them from multiple views. The ultimate output of our method should be a global texture map and material map. Thus, we rasterize the input mesh using the differentiable renderer [23] and optimize its albedo, roughness and metallic parameters via the MSE loss between the Nvdiffrast [23] rendering results and 3D Gaussian rendering results. With adequately optimized implementation, this process takes only about several seconds to bake the multi-view images to untextured 3D model. After the optimization, our input mesh paired with its albedo and material map is capable of performing physically based rendering in new illumination environments.

4. Experiments

4.1. Implementation Details

Dataset We train our TexGaussian model on two publicly-available datasets: ShapeNet [6] and Objaverse [12]. For 3D objects in the ShapeNet dataset, since they only contain albedo maps without PBR materials, our model is trained to generate RGB textures only. We train our model on four categories of ShapeNet: *bench*, *car*, *chair* and *table* which is consistent with previous works [25, 53]. We curate a subset of Objaverse models encoded with PBR materials and convert those using the specular-glossiness workflow to the metallic-roughness workflow to achieve a consistent PBR representation. This process resulted in a total of 29,200 models. For each 3D object, we render its RGBA image, roughness map, and metallic map from 64 views of 512^2 for training.

Network Architecture Our 3D U-Net is built upon octree-based convolutional neural networks [48] which only operates on non-empty octree leaf nodes. The depth of our constructed octree is set to 8 (resolution 256^2). Our 3D U-Net consists of 5 down-sampling and up-sampling blocks.



Figure 3. Unconditional RGB texture generative results on ShapeNet. Please zoom in for a better inspection of color details.

For the text-conditioned generation, cross-attention layers are only inserted at the last two down-sampling blocks, the middle block, and the first two up-sampling blocks. The input channel of our 3D U-Net is set to 4, where 3 for normal vector and 1 for local displacement while the output channel is set to 13 as mentioned above. The resolution of 3D Gaussian Splatting is set to 512^2 . The resolution of the baked albedo and material map is set to 1024^2 .

Training Details We train our model in a single-category and unconditional manner on ShapeNet, i.e., there are 4 TexGaussian models without text conditions in total. This per-category model is trained on 4 NVIDIA A100 (40G) GPUs for about 2 days and is set up for a fair comparison with other existing methods trained on ShapeNet. For the Objaverse dataset, we train a text-conditioned TexGaussian model for PBR material generation on 24 NVIDIA A100 (40G) GPUs for two weeks. For each batch, we randomly sample 8 camera views to calculate loss functions. We adopt the AdamW [26] optimizer with a learning rate of 4×10^{-4} , a weight decay of 0.05, and betas of (0.9, 0.95). The learning rate is cosine annealed to 0 during the training and we clip the gradient with a maximum norm of 1.0.

Evaluation metrics To effectively assess the quality of our generative results, we adopt the metric proposed by [59]. Specifically, each mesh with the generated PBR material is rendered from 20 uniformly distributed views to get multi-view albedo maps, and PBR rendering results in a new illumination environment. These images are used to calculate the FID [35] and KID [1] scores against those from ground truth to evaluate the quality and diversity of generative PBR materials. For the ShapeNet dataset, we only use albedo maps. The final score is averaged across 20 views, and a lower FID and KID score indicates better generation quality and diversity.

4.2. Unconditional RGB Texture Generation

We conduct unconditional RGB texture instead of PBR material generation on four categories of the ShapeNet dataset. We use the same train and test data split as in [25, 53]. Fig. 3

Table 1. Quantitative comparison of the FID and KID ($\times 10^2$) score as well as the inference time of TexGaussian and other methods on ShapeNet dataset [6]. The top part reports the comparison with methods that also train per-category models on ShapeNet. The bottom part reports the comparison with 2D diffusion-based methods, which select 50 samples from the test set of each category.

Methods	Average		Bench		Car		Chair		Table		Time
	FID↓	KID↓	FID↓	KID↓	FID↓	KID↓	FID↓	KID↓	FID↓	KID↓	
Point-UV (1-Stage) [53]	88.43	5.78	64.96	1.44	186.10	17.57	46.23	1.93	56.42	2.19	39.92s
Point-UV (2-Stage) [53]	61.49	2.67	67.48	1.61	89.38	5.82	41.33	1.53	47.75	1.72	49.75s
TexOct [25]	59.45	2.60	60.46	0.97	90.10	6.11	37.70	1.36	49.52	1.97	17.44s
Ours	49.76	2.07	46.37	0.44	80.20	5.22	29.96	1.09	42.52	1.54	11.02s
TEXTure [40]	169.82	7.94	152.53	5.28	236.58	18.52	119.71	2.86	170.46	5.11	132.21s
Text2Tex [7]	156.62	6.68	154.95	5.79	188.60	15.06	125.34	3.11	157.57	2.77	1005.58s
Paint3D [54]	135.25	4.67	104.16	1.68	204.22	13.06	104.64	2.24	127.98	1.71	231.56s
Ours	100.97	1.88	86.37	0.75	117.53	5.20	94.89	0.57	105.09	0.95	11.12s

shows the unconditionally generated RGB texture by our single-category model with high quality, fidelity, and diversity. It can be seen that the texture map synthesized by our method is of great 3D global consistency with the corresponding input geometry.

Quantitative Comparison We conduct quantitative analysis and comparison on our model and other state-of-the-art methods. Specifically, we compare our model with Point-UV Diffusion [53] and TexOct [25] which also train single-category models on ShapeNet. We also compare with some methods using pre-trained 2D diffusion prior for multi-view images synthesis such as TEXTure [53], Text2Tex [7] Paint3D [54], and TexPainter [55] whose input text prompts are set to “a *” and “*” is the name of corresponding category. Due to the relatively long time consumption, we only use 50 3D objects selected from the test set to evaluate the methods based on 2D diffusion models. We conduct only a qualitative comparison for TexPainter [55] for its lengthy processing time. Table 1 reports the comparison of FID and KID scores as well as average inference time on a single NVIDIA A100 (40G) GPU. It is worth noting that our model only takes about one second to predict the parameters of each 3D Gaussian and the rest of the time is used for texture baking. From Table 1, we have the following observation. First, our method obtains the best performance in terms of FID and KID. For example, our method outperforms TexOct [25], by an average of 9.69 in FID and 0.53 in KID. These improvements indicate that our method excels at generating high-quality textures. Second, TexGaussian achieves the fastest generative speed in all categories, which is much less time-consuming than other methods.

Qualitative Comparison Fig. 4 provides some qualitative results on the same input 3D meshes by different methods. Point-UV [53] and TexOct [25] use colored point

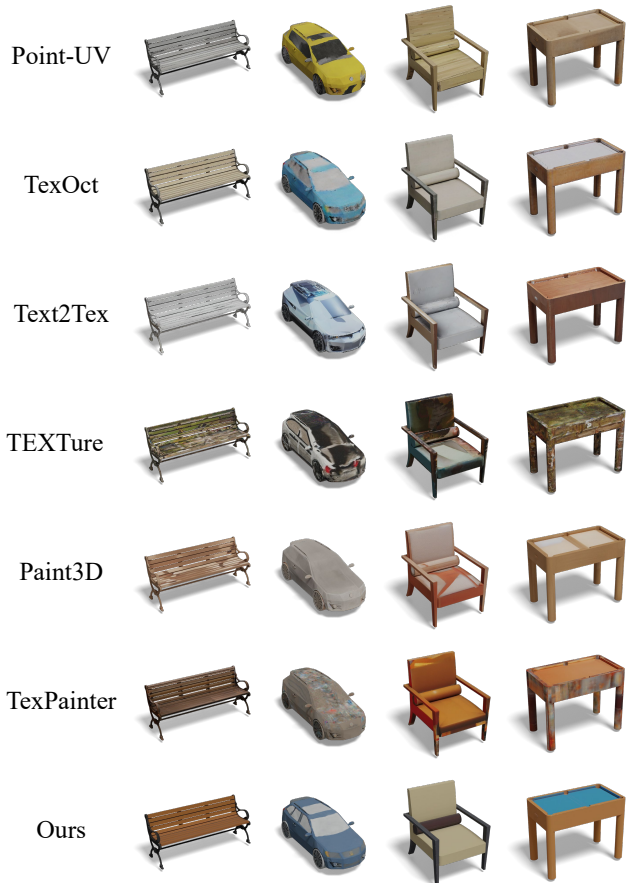


Figure 4. Examples of generated RGB texture obtained by TexGaussian and other state-of-the-art models on the same 3D object. Please zoom in for a better inspection.

clouds as supervision to train diffusion models. As a consequence, the generated texture map is relatively blurry due to the discreteness of the point cloud. For other methods

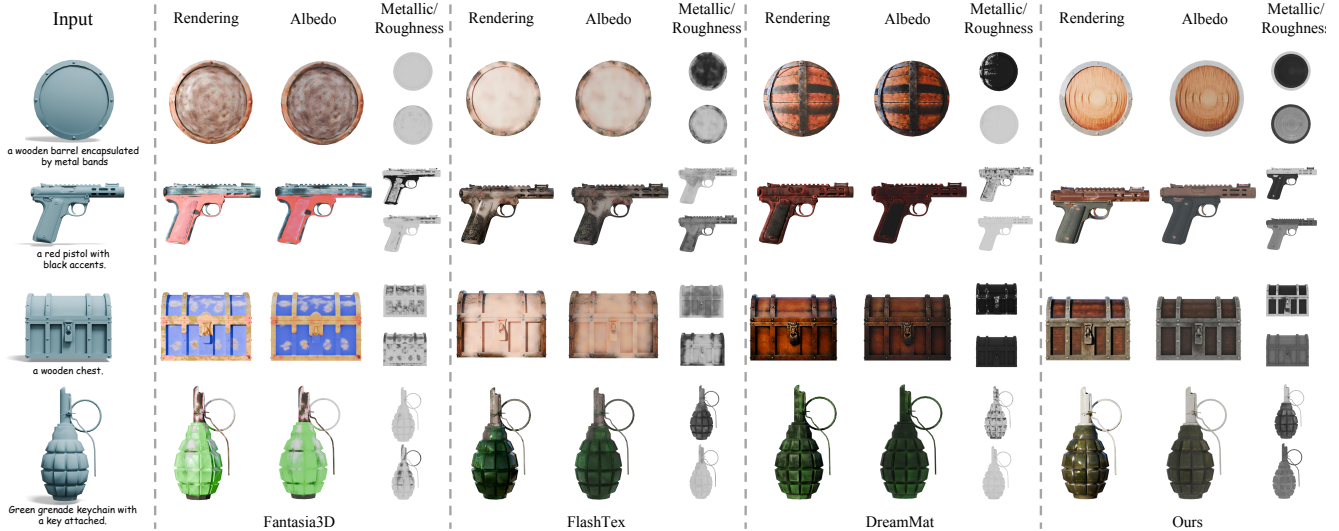


Figure 5. Qualitative comparison with Fantasia3D [8], FlashTex [13] and DreamMat [58]. We provide the rendered image, albedo map, roughness map, and metallic map for each 3D object.

that leverage pre-trained 2D diffusion prior, they can hardly comprehend the total geometry. The generated texture maps by them are not consistent with the semantics of input 3D mesh, such as the chaotic stripes on the bench and chair in Fig. 4. On the contrary, our method generates smooth and colorful RGB textures on unseen objects, which align well with 3D meshes.

4.3. Text-conditioned PBR Material Generation

We train the text-conditioned PBR material generation model on our filtered subset of Objaverse [12]. We use the text descriptions of 3D objects from Cap3D [27, 28] to train our model. We choose 29,000 3D objects in our filtered subset for training and the rest 200 for testing. Fig 1 shows some generative results by proposed TexGaussian on the test set. We can see that our model comprehends the overall geometry feature of input 3D mesh. It is capable of synthesizing high-quality albedo, roughness, and metallic, which are aligned well with 3D semantics such as accents with high metallic on the wooden box and rope with high roughness on the jug.

Quantitative Comparison We conduct quantitative comparison with three state-of-the-art text-conditioned PBR material synthesis methods: Fantasia3D [8], FlashTex [13], and DreamMat [58] on our test set. Table 2 reports FID and KID scores on both multi-view albedo map and PBR rendering images under the same illumination environment. TexGaussian outperforms all the baseline methods in achieving the best visual quality of the generated appearances. We also report the average inference time of different methods on a single NVIDIA A100 (40G) GPU across

Table 2. Quantitative comparison of the FID and KID ($\times 10^2$) scores as well as the inference time of TexGaussian and other methods on our test set which consists of 200 3D objects with ground truth PBR materials.

Methods	Albedo		PBR rendering		Time
	FID↓	KID↓	FID↓	KID↓	
Fantasia3D [8]	213.21	0.96	209.87	0.48	22.7mins
FlashTex [13]	185.24	1.13	186.82	0.42	20.3mins
DreamMat [58]	152.63	1.09	145.49	0.19	48.2mins
Ours	123.72	0.20	129.52	0.02	21.04s

our test set. Due to the iterative optimization of score distillation sampling, all other three methods cost at least 20 minutes for the generation. In contrast, our TexGaussian only takes about 20 seconds, in which one second is for predicting the Gaussian parameters and the rest is for baking, which results in $60\times$ faster compared to previous approaches.

Qualitative Comparison Fig 5 visualizes the generated albedo, roughness, and metallic of each compared method from the same text prompt and untextured meshes. We also show the PBR rendering images of the generated materials under the same environment light. It can be observed that Fantasia3D [8] and FlashTex [13] generate irregular colors on mesh surfaces, which lead to rendering results with low quality. DreamMat [58] is capable of generating a visually pleasing appearance for test 3D models. However, it can hardly align the generated PBR material well with the 3D semantics, such as the black band on the wooden barrel. What’s more, it tends to generate over-saturated colors, which is demonstrated in the wooden chest and green

Table 3. Quantitative comparison of methods using different types of 3D Gaussian Fitting.

Method	PSNR \uparrow	LPIPS \downarrow	SSIM \uparrow	3D Gaussian Number
Full voxel aligned	23.75	0.17	0.92	16,777,216 (256^3)
Octant-aligned	32.60	0.039	0.97	79,480

grenade. In addition, all of the compared methods struggle to completely disentangle the light and texture, which results in relatively chaotic results of generated metallic and roughness. On the contrary, our method is directly trained from 3D origin data. It is capable of generating clean and smooth PBR material while fully comprehending the overall 3D structure. We provide more generative results in the supplementary material.

4.4. Ablation Study

For the purpose of analyzing the impact of different designs in our model, we conduct ablation studies by removing or changing some proposed modules.

Essentials of Octree We first analyze the effectiveness of using Octree in our pipeline. To do this, we calculate the quantitative quality of 3D Gaussian pre-fitting by our octant-aligned and full voxel-aligned 3D Gaussian Splatting on the subset of Objaverse dataset [12]. The depth of octree in this ablation is set to 8, and the resolution of full voxels is 256^3 , which is consistent with the finest resolution of octree to guarantee comparison fairness. Table 3 reports the PSNR, LPIPS [57], and SSIM [49] metrics on the albedo map as well as the average number of 3D Gaussians of octant and full voxel aligned 3D Gaussian Splatting across our Objaverse subset. Fig. 6 also presents some visual results of our Gaussian fitting. These results demonstrate that the proposed octant-aligned 3D Gaussian produces a much more reasonably compact and precise representation of diverse and complex 3D assets with much fewer 3D Gaussians compared to the full voxel version.

Essentials of different losses We verify the importance of different losses we proposed to train our model. We term all the losses calculated on 2D image space as L_{2D} :

$$L_{2D} = L_{RGB} + L_R + L_M + L_\alpha. \quad (7)$$

We train two additional TexGaussian models using only the L_{2D} or L_{3D} loss on ShapeNet `car` category due to its large variations and complexity of texture to validate their effectiveness. It is worth noting that in this dataset, $L_{2D} = L_{RGB} + L_\alpha$ due to the lack of material information. The training curves of MSE loss and LPIPS loss between rendering images and ground-truth ones are shown



Figure 6. Visualization of different manners of Gaussian fitting. The rendering results demonstrate excellent reconstruction performance of the proposed octant-aligned 3D Gaussian Splatting.

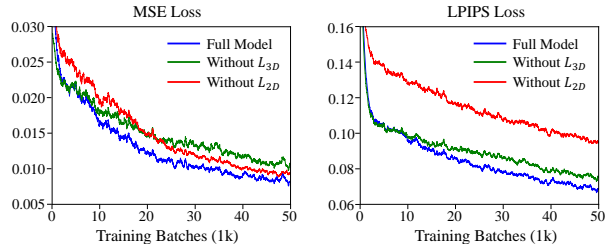


Figure 7. The training loss curves of our model with different losses.

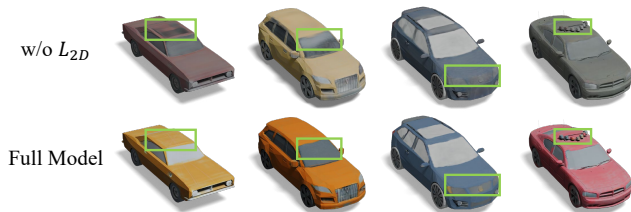


Figure 8. Results of our model trained with and without L_{2D} . The green rectangles highlight the shortcomings of TexGaussian model trained without L_{2D} .

in Fig 7, which demonstrate the effects of L_{2D} and L_{3D} . From the loss curves, we can conclude that L_{3D} facilitates the process of convergence and L_{2D} enhances the quality of synthesized images, as evidenced by the large margin of improvement in the LPIPS loss when introducing L_{2D} . We also provide some qualitative results to verify the effectiveness of L_{2D} on the test set of `car` category in Fig 8. Only using L_{3D} results in a relatively blurry texture map due to the discreteness of 3D Gaussian which is similar to the characteristic of 3D point clouds analyzed above.

5. Conclusion

In this paper, we proposed TexGaussian, an octree-based 3D Gaussian Splatting model for high-quality PBR material generation on untextured meshes. We aligned each 3D

Gaussian on the octant of the corresponding octree built from the input untextured object and extended the parameters of 3D Gaussian with additional channels to represent the roughness and metallic map. We trained our model with regression objectives, achieving faster inference speed compared to previous texture synthesis methods. Experimental results demonstrated that our method is capable of generating high-quality PBR materials that are readily usable in modern graphics engines for photo-realistic rendering, offering enhanced realism for a variety of applications.

Limitations The generalization of TexGaussian is still hindered by the scale of the training set. Thus it struggles to generate various textures for some extremely complex 3D objects beyond our training data. We are looking forward to training our TexGaussian model with more parameters and more data on a larger-scale GPU cluster in the future.

References

- [1] Mikołaj Bińkowski, Dougal J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying MMD GANs. In *International Conference on Learning Representations*, 2018. 5
- [2] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 2
- [3] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22563–22575, 2023.
- [4] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. 2024. 2
- [5] Tianshi Cao, Karsten Kreis, Sanja Fidler, Nicholas Sharp, and KangXue Yin. Textfusion: Synthesizing 3d textures with text-guided image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 2, 3
- [6] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical Report arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago, 2015. 3, 5, 6
- [7] Dave Zhenyu Chen, Yawar Siddiqui, Hsin-Ying Lee, Sergey Tulyakov, and Matthias Nießner. Text2tex: Text-driven texture synthesis via diffusion models. *arXiv preprint arXiv:2303.11396*, 2023. 2, 3, 6
- [8] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 22246–22256, 2023. 3, 7
- [9] Yongwei Chen, Rui Chen, Jiabao Lei, Yabin Zhang, and Kui Jia. Tango: Text-driven photorealistic and robust 3d stylization via lighting decomposition. *Advances in Neural Information Processing Systems*, 35:30923–30936, 2022. 3
- [10] Jasmine Collins, Shubham Goel, Kenan Deng, Achleshwar Luthra, Leon Xu, Erhan Gundogdu, Xi Zhang, Tomas F Yago Vicente, Thomas Dideriksen, Himanshu Arora, et al. Abo: Dataset and benchmarks for real-world 3d object understanding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 21126–21136, 2022. 3
- [11] 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.
- [12] 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. 2, 3, 5, 7, 8
- [13] Kangle Deng, Timothy Omerick, Alexander Weiss, Deva Ramanan, Jun-Yan Zhu, Tinghui Zhou, and Maneesh Agrawala. Flashtex: Fast relightable mesh texturing with lightcontrolnet. In *European Conference on Computer Vision (ECCV)*, 2024. 3, 7
- [14] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 2
- [15] Chenjian Gao, Boyan Jiang, Xinghui Li, Yingpeng Zhang, and Qian Yu. Genesisstex: Adapting image denoising diffusion to texture space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4620–4629, 2024. 3
- [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. 2020. 2
- [17] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022. 2
- [18] Fangzhou Hong, Mingyuan Zhang, Liang Pan, Zhongang Cai, Lei Yang, and Ziwei Liu. Avatarclip: Zero-shot text-driven generation and animation of 3d avatars. *arXiv preprint arXiv:2205.08535*, 2022. 3
- [19] Dong Huo, Zixin Guo, Xinxin Zuo, Zhihao Shi, Juwei Lu, Peng Dai, Songcen Xu, Li Cheng, and Yee-Hong Yang. Texgen: Text-guided 3d texture generation with multi-view sampling and resampling. *ECCV*, 2024. 2, 3
- [20] Yingwenqi Jiang, Jiadong Tu, Yuan Liu, Xifeng Gao, Xiaoxiao Long, Wenping Wang, and Yuexin Ma. Gaussian-

- shader: 3d gaussian splatting with shading functions for reflective surfaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5322–5332, 2024. 4
- [21] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), 2023. 2, 4
- [22] Matthias Labschütz, Katharina Krösl, Mariebeth Aquino, Florian Grashärtl, and Stephanie Kohl. Content creation for a 3d game with maya and unity 3d. *Institute of Computer Graphics and Algorithms, Vienna University of Technology*, 6(124):2, 2011. 1
- [23] Samuli Laine, Janne Hellsten, Tero Karras, Yeongho Seol, Jaakko Lehtinen, and Timo Aila. Modular primitives for high-performance differentiable rendering. *ACM Transactions on Graphics*, 39(6), 2020. 5
- [24] 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 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3
- [25] Jialun Liu, Chenming Wu, Xinqi Liu, Xing Liu, Jinbo Wu, Haotian Peng, Chen Zhao, Haocheng Feng, Jingtuo Liu, and Errui Ding. Textoc: Generating textures of 3d models with octree-based diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4284–4293, 2024. 2, 3, 5, 6
- [26] I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 5
- [27] Tiange Luo, Chris Rockwell, Honglak Lee, and Justin Johnson. Scalable 3d captioning with pretrained models. *arXiv preprint arXiv:2306.07279*, 2023. 7
- [28] Tiange Luo, Justin Johnson, and Honglak Lee. View selection for 3d captioning via diffusion ranking. *arXiv preprint arXiv:2404.07984*, 2024. 7
- [29] Yiwei Ma, Xiaoqing Zhang, Xiaoshuai Sun, Jiayi Ji, Haowei Wang, Guannan Jiang, Weilin Zhuang, and Rongrong Ji. X-mesh: Towards fast and accurate text-driven 3d stylization via dynamic textual guidance. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2749–2760, 2023. 3
- [30] Oscar Michel, Roi Bar-On, Richard Liu, Sagie Benaim, and Rana Hanocka. Text2mesh: Text-driven neural stylization for meshes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13492–13502, 2022. 3
- [31] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 3
- [32] Nasir Mohammad Khalid, Tianhao Xie, Eugene Belilovsky, and Tiberiu Popa. Clip-mesh: Generating textured meshes from text using pretrained image-text models. In *SIGGRAPH Asia 2022 conference papers*, pages 1–8, 2022. 3
- [33] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, 2022. 3
- [34] Michael Oechsle, Lars Mescheder, Michael Niemeyer, Thilo Strauss, and Andreas Geiger. Texture fields: Learning texture representations in function space. In *Proceedings IEEE International Conf. on Computer Vision (ICCV)*, 2019. 3
- [35] Gaurav Parmar, Richard Zhang, and Jun-Yan Zhu. On aliased resizing and surprising subtleties in gan evaluation. In *CVPR*, 2022. 5
- [36] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image synthesis. In *The Twelfth International Conference on Learning Representations*, 2024. 2
- [37] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 2, 3
- [38] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 3, 4
- [39] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022. 2
- [40] Elad Richardson, Gal Metzer, Yuval Alaluf, Raja Giryes, and Daniel Cohen-Or. Texture: Text-guided texturing of 3d shapes. In *ACM SIGGRAPH 2023 conference proceedings*, pages 1–11, 2023. 2, 3, 6
- [41] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021. 2
- [42] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022. 2
- [43] Yahao Shi, Yanmin Wu, Chenming Wu, Xing Liu, Chen Zhao, Haocheng Feng, Jingtuo Liu, Liangjun Zhang, Jian Zhang, Bin Zhou, et al. Gir: 3d gaussian inverse rendering for relightable scene factorization. *arXiv preprint arXiv:2312.05133*, 2023. 4
- [44] Yawar Siddiqui, Justus Thies, Fangchang Ma, Qi Shan, Matthias Nießner, and Angela Dai. Texturify: Generating textures on 3d shape surfaces. In *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part III*, pages 72–88. Springer, 2022. 3
- [45] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015. 2
- [46] Jiayang 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. 4
- [47] Peng-Shuai Wang. Octformer: Octree-based transformers for 3D point clouds. *ACM Transactions on Graphics (SIGGRAPH)*, 42(4), 2023. 4
- [48] Peng-Shuai Wang, Yang Liu, Yu-Xiao Guo, Chun-Yu Sun, and Xin Tong. O-CNN: Octree-based convolutional neural networks for 3D shape analysis. *ACM Transactions on Graphics (SIGGRAPH)*, 36(4), 2017. 2, 4, 5
- [49] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 8
- [50] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 3
- [51] Jinbo Xing, Hanyuan Liu, Menghan Xia, Yong Zhang, Xintao Wang, Ying Shan, and Tien-Tsin Wong. Toon-crafter: Generative cartoon interpolation. *arXiv preprint arXiv:2405.17933*, 2024. 2
- [52] Yu-Ying Yeh, Jia-Bin Huang, Changil Kim, Lei Xiao, Thu Nguyen-Phuoc, Numair Khan, Cheng Zhang, Manmohan Chandraker, Carl S Marshall, Zhao Dong, et al. Texture-dreamer: Image-guided texture synthesis through geometry-aware diffusion. *arXiv preprint arXiv:2401.09416*, 2024. 3
- [53] Xin Yu, Peng Dai, Wenbo Li, Lan Ma, Zhengzhe Liu, and Xiaojuan Qi. Texture generation on 3d meshes with point-uv diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4206–4216, 2023. 2, 3, 5, 6
- [54] Xianfang Zeng, Xin Chen, Zhongqi Qi, Wen Liu, Zibo Zhao, Zhibin Wang, Bin Fu, Yong Liu, and Gang Yu. Paint3d: Paint anything 3d with lighting-less texture diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4252–4262, 2024. 2, 3, 6
- [55] Hongkun Zhang, Zherong Pan, Congyi Zhang, Lifeng Zhu, and Xifeng Gao. Texpainter: Generative mesh texturing with multi-view consistency. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–11, 2024. 2, 3, 6
- [56] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023. 2, 3
- [57] 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. 4, 8
- [58] Yuqing Zhang, Yuan Liu, Zhiyu Xie, Lei Yang, Zhongyuan Liu, Mengzhou Yang, Runze Zhang, Qilong Kou, Cheng Lin, Wenping Wang, and Xiaogang Jin. Dreammat: High-quality pbr material generation with geometry- and light-aware diffusion models. *ACM Trans. Graph.*, 43(4), 2024. 2, 3, 7
- [59] Xin-Yang Zheng, Hao Pan, Peng-Shuai Wang, Xin Tong, Yang Liu, and Heung-Yeung Shum. Locally attentional sdf diffusion for controllable 3d shape generation. *ACM Transactions on Graphics (SIGGRAPH)*, 42(4), 2023. 5

A. Preliminary of 3D Gaussian Splatting

Gaussian splatting employs a collection of 3D Gaussians to represent 3D data. Specifically, each Gaussian is formally defined as:

$$G(\mathbf{x}) = e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}, \quad (8)$$

where $\boldsymbol{\mu} \in \mathbb{R}^3$ represents the spatial mean of 3D Gaussian and $\boldsymbol{\Sigma} \in \mathbb{R}^{3 \times 3}$ denotes the covariance matrix. The covariance matrix $\boldsymbol{\Sigma}$ of a 3D Gaussian is analogous to describing the configuration of an ellipsoid. Thus, the covariance matrix $\boldsymbol{\Sigma}$ is decomposed into a scaling matrix \mathbf{S} and a rotation matrix \mathbf{R} as follows:

$$\boldsymbol{\Sigma} = \mathbf{R} \mathbf{S} \mathbf{S}^T \mathbf{R}^T, \quad (9)$$

To allow independent optimization of both factors, they are stored separately: a 3D vector \mathbf{s} for scaling and a quaternion \mathbf{q} to represent rotation. During the rendering process, the 3D Gaussians are projected onto a 2D plane. With the intrinsic matrix \mathbf{K} and extrinsic matrix \mathbf{T} , the 2D mean $\boldsymbol{\mu}'$ and covariance $\boldsymbol{\Sigma}'$ are defined as follows:

$$\boldsymbol{\mu}' = \mathbf{K}[\boldsymbol{\mu}, 1]^T, \quad \boldsymbol{\Sigma}' = \mathbf{J} \mathbf{T} \boldsymbol{\Sigma} \mathbf{T}^T \mathbf{J}^T, \quad (10)$$

Here, \mathbf{J} represents the Jacobian of the affine approximation of the projective transformation. Each 3D Gaussian is associated with an opacity value o and a view-dependent color \mathbf{c} , determined by a set of spherical harmonics coefficients. In our model, the multi-view rendered images of albedo map do not depend on the selected viewpoints. As a result, we just use three-channels RGB on each 3D Gaussian to represent the view-independent colors instead of original spherical harmonics, and we exclude the positional parameter $\boldsymbol{\mu}$ because each 3D Gaussian is fixed at the center of each finest leaf node of the constructed octree. All the parameters can be collectively denoted by Θ_0 with:

$$\Theta_{0_i} = \{\mathbf{o}_i, \mathbf{s}_i, \mathbf{q}_i, \mathbf{c}_i\}, \quad (11)$$

representing the parameters for the i -th Gaussian.

Moreover, to encode the PBR material parameters, we append additional two parameters: roughness r and metallic m at the end of the original Gaussian parameters. To render multi-view images of these two attributes, we concatenate r and m with previous parameters to obtain:

$$\Theta_{1_i} = \{\mathbf{o}_i, \mathbf{s}_i, \mathbf{q}_i, \mathbf{r}_i\}, \quad \Theta_{2_i} = \{\mathbf{o}_i, \mathbf{s}_i, \mathbf{q}_i, \mathbf{m}_i\}. \quad (12)$$

Then, all the 3D Gaussians are paired with these two new parameters Θ_{1_i} and Θ_{2_i} , rendered from multiple viewpoints to get multi-view roughness map and metallic map for further training.

B. Network Details

Unconditional RGB Texture Generation The network architecture of the octree-based 3D U-Net we used in unconditional RGB texture generation is shown in Fig 9. The U-Net has five hierarchical levels, corresponding to octree depths of 8, 7, 6, 5 and 4, with resolutions of 256^3 , 128^3 , 64^3 , 32^3 , 16^3 . The feature dimensions are set to 32, 64, 128, 256, 256 respectively. The channel of input and output feature is 4 and 13 as described in the main manuscript.

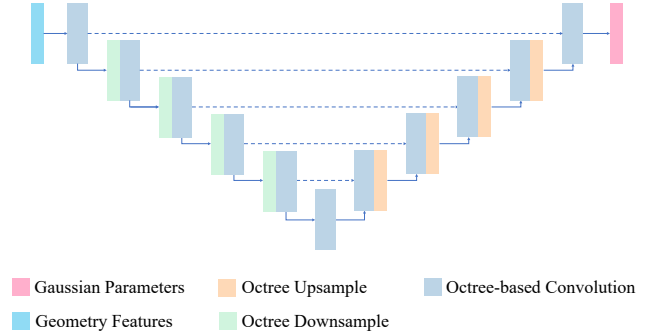


Figure 9. The network architecture of the octree-based 3D U-Net we used to train our unconditional RGB texture generation model.

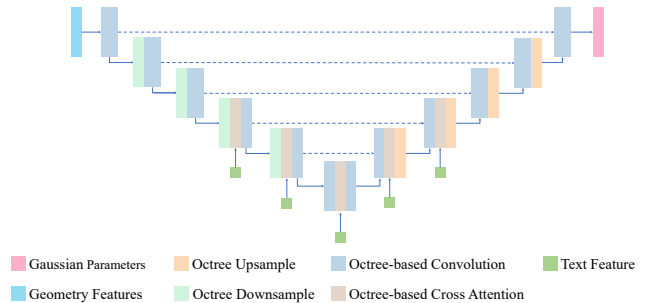


Figure 10. The network architecture of the octree-based 3D U-Net we used to train our text-conditioned PBR material generation model.

Text-conditioned PBR Material Generation The network architecture of the octree-based 3D U-Net we used in text-conditioned PBR material generation is shown in Fig 10. The U-Net has five hierarchical levels, corresponding to octree depths of 8, 7, 6, 5 and 4, with resolutions of 256^3 , 128^3 , 64^3 , 32^3 , 16^3 . The feature dimensions are set to 64, 128, 256, 512, 512 respectively. The text feature is fed to U-Net via the octree-based multi-head cross attention mechanism. The cross attention layers are only inserted at the least two down-sampling blocks, the middle block and the two first up-sampling blocks to save GPU memory.

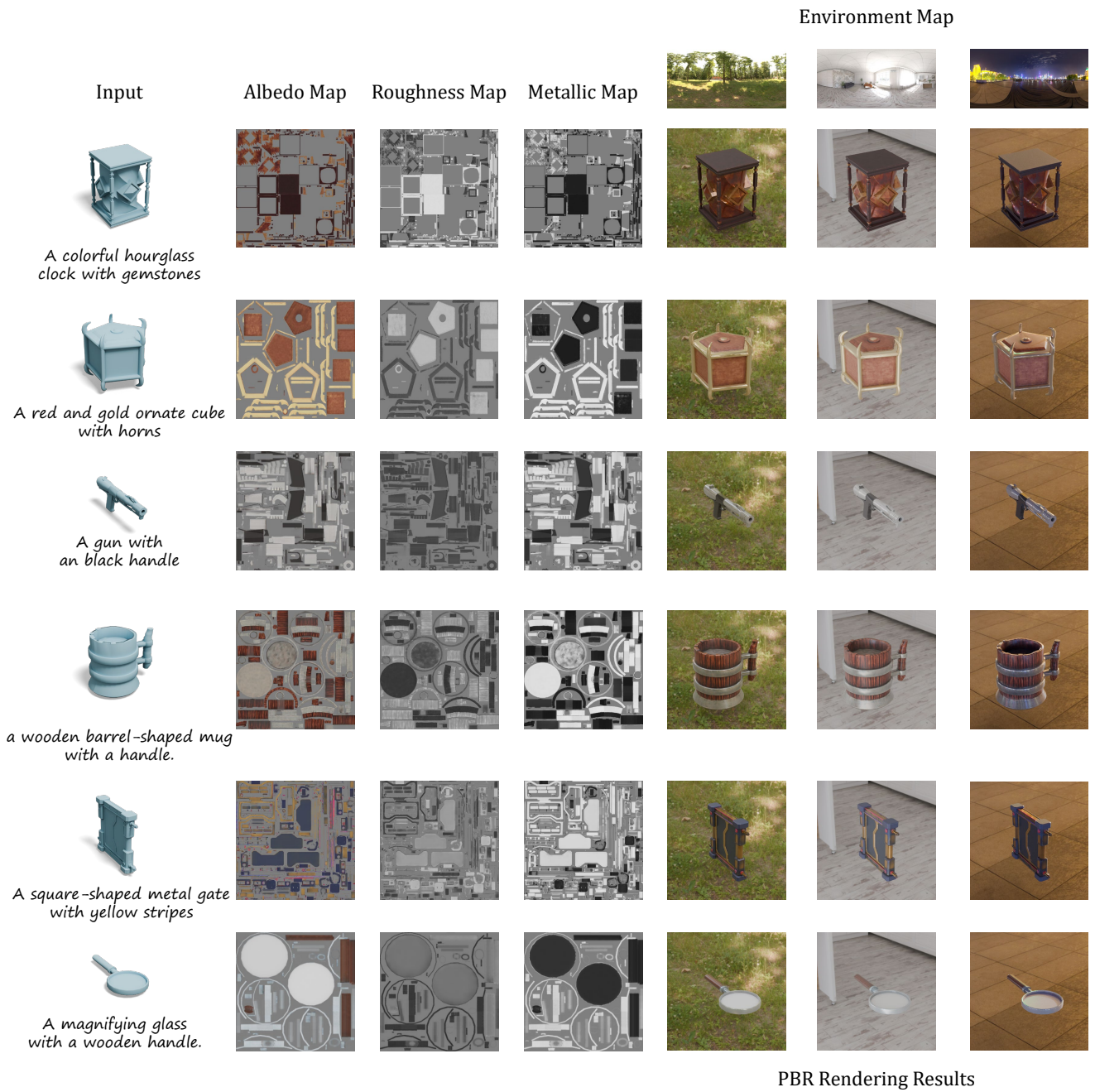
C. More Results

Our method is capable of generating diverse materials given different text prompts for a single mesh. Fig 11 shows the PBR materials and the rendering results of the same mesh generated from different text prompts by our proposed TexGaussian. These results demonstrate that our method is able to generate diverse materials of different styles that align well with the text prompts and 3D objects with high fidelity.

We provide more generated results in Fig 12.



Figure 11. Diverse material generation. Our method can generate different materials with different text prompts on the same mesh.



PBR Rendering Results

Figure 12. More generative results of our method on different input 3D models and text prompts.