
Decorate3D: Text-Driven High-Quality Texture Generation for Mesh Decoration in the Wild

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 This paper introduces Decorate3D, a flexible and easy-to-use tool to create & edit
2 3D objects from images. Given multi-view photos of an object of interest, our
3 approach starts with a textured mesh 3D reconstruction of the object, followed
4 by either manual editing, or text-to-3D texture automatic generation of textures
5 that align with the mesh. Decorate3D decomposes the Neural Radiance Field
6 into an explicit mesh representation and a view-dependent texture, allowing for
7 flexible decoration. One key contribution of Decorate3D is paving the way to
8 text-driven 3D high-quality texture generation, which has remained an unresolved
9 problem. To tackle this challenge, we propose a structure-aware score distillation
10 sampling method to optimize a UV neural texture based on user-defined text
11 and empower an image diffusion model with 3D-consistent generation capability.
12 Furthermore, we introduce a few-view resampling training method and utilize a
13 super-resolution model to obtain high-quality UV textures (2048×2048) for 3D
14 texturing. Extensive experiments collectively validate the superior performance of
15 Decorate3D in decorating real-world 3D objects.

1 Introduction

17 The recent development of effective neural 3D reconstruction techniques such as neural radiance field
18 (NeRF) [11] has facilitated the creation of realistic digital replicas for real-world 3D scenes. While
19 significant progress has been made in realistic 3D reconstruction, one remaining challenge is how
20 to allow users to edit or decorate the acquired 3D objects or assets in the digital scene. To address
21 this issue, this paper introduces Decorate 3D, a user-friendly approach to editing 3D objects through
22 either easy-to-use manual editing or text prompt guided texture generation, as showcased in Fig. 1.

23 Since the implicit representations of the NeRF model are tightly coupled, it is not trivial to achieve
24 the aforementioned operations for the 3D scene’s decoration. To facilitate the modifiability of the 3D
25 representations, we bake the NeRF model into a triangle mesh representation with a view-dependent
26 texture in UV space, which is called a *Decomposition* phase in this paper. The decoupled geometry
27 and UV texture representation make the subsequent *Decoration* phase more manageable. In the
28 decoration phase, Decorate3D allows flexible texture editing and controllable generation of UV
29 textures with the instruction of prompts.

30 The success of the decoration phase demands a solution to the text-to-3D synthesis problem, which
31 has been so far restricted by the lack of paired text and 3D assets. Recent advancements in text-
32 to-image models [22, 18, 14] greatly facilitate text-driven 3D editing or generation. [31, 1, 15, 10].
33 Researchers have proposed approaches using pre-trained text-to-2D diffusion models to optimize
34 NeRF in the zero-shot setting via the Score Distillation Sampling (SDS) [15] strategy. Although
35 impressive results have been achieved, the naive SDS optimization based on 2D diffusion models
36 lacks 3D awareness, resulting in inharmonious textures that misalign with the geometry.

37 To address the above issues, during the decoration phase of Decorate3D , we carry out a direct update
 38 of the global UV texture map via a structure-aware SDS optimization. We employ a pre-trained
 39 depth-guided text-to-image latent diffusion model [28, 22] and incorporate the initial decoupled UV
 40 texture from the decomposition phase as an additional form of structural guidance. We empirically
 41 find that rendering with the straightforwardly optimized UV texture from SDS tends to be noisy
 42 with color distortions. The reason is that the optimized UV texture stands for neural features in
 43 effect, which produce rendered neural images that necessitate a neural interpreter. We incorporate the
 44 idea of deferred neural rendering [30] and synthesize images of different views by forwarding the
 45 rendered neural images to the encoder-decoder of the latent diffusion model. However, the decoder’s
 46 view-by-view synthesis causes jittering effects. To address this issue, a few-view resampling training
 47 strategy is proposed to optimize a global UV texture from sparsely sampled view directions. Finally,
 48 to further enhance the visual quality, we use a super-resolution diffusion model to improve the
 49 resolution, which is directly applied to the UV space.

50 Our key contributions can be summarized as follows. First, we propose a system pipeline to allow
 51 convenient 3D-consistent decoration of 3D objects captured in the wild; Second, with our structure-
 52 aware 3D texture generation, few-view resampling training, and super-resolution enhancement, we
 53 are able to synthesize high-quality textures aligned well with the geometry; Finally, we demonstrate
 54 the effectiveness of Decorate3D on real-world datasets and have achieved superior performance over
 55 state-of-the-art approaches.



Figure 1: Given captured images of an object, **Decorate3D** supports text-driven high-quality texture generation and user-friendly texture editing. Please zoom in for better visualization.

56 2 Related Work

57 **Text-to-Image Diffusion Models** The past two years have witnessed the success of multiple large
 58 diffusion models [19, 5, 22, 18, 14] that are able to generate impressive images with photo-realistic
 59 details conditioned on an input text prompt. The widely popular stable diffusion [22], was trained
 60 on rich paired text-image data and is conditioned on CLIP’s [16] frozen text encoder. Beyond
 61 text-conditioning, ControlNet [35] extended the stable diffusion by training a parallel hyper-network
 62 that allows controllable generation with additional input modalities such as depth maps or edges.

63 **Text-to-3D Generation Methods** The development of 2D image generation also greatly facilitates
 64 the techniques of text-to-3D generation. CLIP-Mesh [7] proposed a text-driven 3D content generation
 65 method using a pre-trained CLIP model. DreamFusion [15] first proposed a score distillation sampling
 66 (SDS) method to achieve text-driven optimization for NeRF, relying on a text-to-image diffusion
 67 model. Magic3D [9] improved DreamFusion’s resolution using a coarse-to-fine optimization strategy.
 68 Latent-Paint [10] proposed to bring the NeRF to the latent space and apply the SDS to optimize a
 69 latent code for 3D scene generation. More recently, Fantasia3d [3] used a hybrid scene representation
 70 of DMTET for SDS optimization, and 3DFuse [26] improved the 3D consistency by incorporating a
 71 coarse 3D prior into fine-tuning the diffusion model. Shape-E [6] trained a conditional 3D diffusion
 72 model using paired 3D and text data. Concurrent works, TEXTure [21] as well as Text2Tex [2],
 73 introduced optimization-free iterative update schemes to paint 3D models from different viewpoints
 74 using depth-to-image diffusion models. Although these optimization-free methods may generate
 75 plausible results for simple geometries, they usually fail on complex surfaces, resulting in artifacts on
 76 the seamed areas. Decorate3D distinguishes itself from previous methods by providing 3D-consistent
 77 and high-quality structure-aware texture generation for diverse objects.

78 3 Approach

79 Decorate3D is a two-phase framework to enable 3D-consistent editing of real-world objects. The
 80 overall framework of Decorate3D is illustrated in Fig. 2. It consists of a decomposition phase and

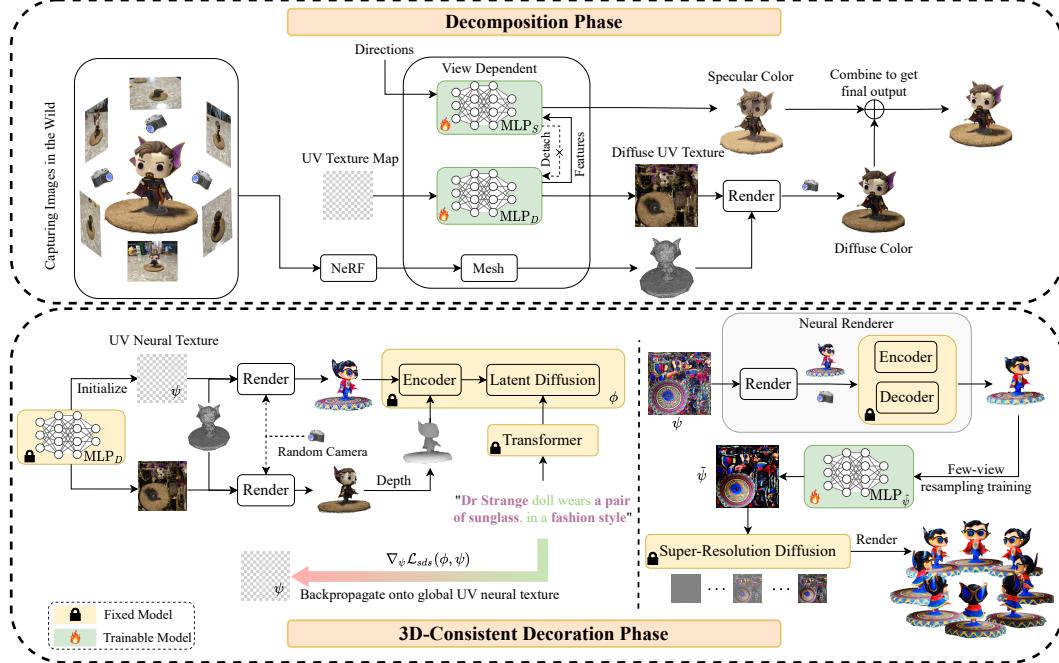


Figure 2: Overview of **Decorate3D**: In the decomposition phase, we utilize NeuS [32] to extract a 3D mesh, and optimize two MLPs (*i.e.* MLP_D and MLP_S) to model view-dependent textures. In the decoration phase, we initialize the UV neural texture ψ using MLP_D , and optimize the neural texture ψ via the depth-guided diffusion model with the score distillation sampling loss. The neural texture is rendered into sampled views and forwarded to the Neural Renderer to synthesize color images. Afterward, we optimize a MLP through the few-view resampling training method to obtain a UV texture $\tilde{\psi}$ in RGB space. Finally, we enhance the visual quality of $\tilde{\psi}$ with a super-resolution diffusion model.

81 a decoration phase. Next, we will introduce the decomposition phase in Sec 3.1 followed by the
 82 decoration phase. The decoration phase is divided into three parts, including text-driven UV neural
 83 texture optimization (Sec 3.2), few-view resampling training (Sec 3.3), and super-resolution on UV
 84 texture (Sec 3.4).

85 3.1 Decomposition Phase

86 Given the captured multi-view images of an object and their respective camera poses obtained by
 87 COLMAP [25], we can easily train a NeRF model for novel view synthesis. But NeRF editing
 88 poses challenges due to its tightly coupled representations. To address these challenges and facilitate
 89 convenient 3D editing, as well as seamless integration into downstream applications, we propose
 90 a solution: decomposing the NeRF representations into a 3D mesh and a view-dependent texture.
 91 Specifically, we employ NeuS [32] to reconstruct a triangle mesh \mathcal{M} of the object. And we calculate
 92 the UV atlas associated with the mesh using XAtlas [34]. To represent the UV texture, we employ
 93 MLP_D to represent a global UV feature map and a diffuse texture, and MLP_S conditioned on view
 94 direction to represent the view-dependent specular texture.

95 Denoting the differentiable mesh rendering by \mathcal{R} , we formulate the image generation process as
 96 follows:

$$\mathbf{c}_d, \mathbf{f}_s = \text{sigmoid}(\text{MLP}_D(\mathbf{v})) \quad (1)$$

$$\mathbf{c}_s = \text{MLP}_S(\mathbf{f}_s, \mathbf{d}) \quad (2)$$

$$\mathbf{c} = \mathcal{R}(\mathbf{c}_d + \mathbf{c}_s, \mathcal{M}) \quad (3)$$

97 where \mathbf{v} refers to 2D UV coordinates after applying positional encoding, and \mathbf{f}_s are the intermediate
 98 features to synthesize the specular color. The features \mathbf{f}_s are forwarded to the view-dependent network
 99 MLP_S to generate view-dependent effects, conditioned on the view direction \mathbf{d} . The final rendered
 100 image \mathbf{c} is obtained by summing up the diffuse color \mathbf{c}_d and the specular color \mathbf{c}_s .

101 To train the texture networks, we first adopt a typical reconstruction loss \mathcal{L}_{color} to minimize the
 102 difference between the rendered image $\hat{\mathbf{c}}(\mathbf{x})$ and its corresponding captured ground truth image $\mathbf{c}(\mathbf{x})$
 103 at each pixel \mathbf{x} . The loss function is formulated as follows:

$$\mathcal{L}_{color} = \sum_{\mathbf{x}} \|\hat{\mathbf{c}}(\mathbf{x}) - \mathbf{c}(\mathbf{x})\|^2 \quad (4)$$

104 In addition, we apply an L1 regularization to enforce sparsity on the specular color:

$$\mathcal{L}_{specular} = \sum_i |\mathbf{c}_s(\mathbf{x}_i)| \quad (5)$$

105 The training of MLP_D and MLP_S can be achieved with a single-stage framework, as illustrated in the
 106 upper part of Fig. 2. More specifically, MLP_D and MLP_S are jointly optimized, where the gradient
 107 from the specular branch will only be used to update MLP_S , and not be propagated to MLP_D . The
 108 MLP_S is a by-product for decomposing the diffuse and specular textures, not used for the decoration
 109 phase. We follow instant-npg [13] to accelerate the optimization to achieve convergence within 3
 110 minutes on a single NVIDIA V100 GPU. More details can be found in the supplementary.

111 3.2 Text-Driven Neural Texture Optimization

112 The text-driven texture generation is the core element of the decoration phase of Decorate3D .
 113 However, diversified 3D generation is often infeasible due to a lack of enough data pairs of text and
 114 3D models. In this context, we utilize the SDS technique for the optimization of the neural texture.

115 **Preliminary.** Let us first introduce the SDS proposed by DreamFusion [15], which achieved text-
 116 to-3D generation based on a pre-trained text-to-image diffusion-based generative model [23]. They
 117 utilized a NeRF representation to model the scene. It is a parametric function $\mathbf{x} = \Gamma(\xi)$, which
 118 can synthesize an image \mathbf{x} at the desired camera pose. Here, Γ is a volumetric renderer, and ξ is a
 119 MLP representing a NeRF scene. The diffusion model ϕ contains a denoising function $\epsilon_\phi(\mathbf{x}_t; \mathbf{y}, t)$
 120 that predicts the sampled noise given the noisy image \mathbf{x}_t at timestep t , and text embedding \mathbf{y} . The
 121 difference between the noisy image and the denoised image provides the gradient to update ξ such
 122 that the density regions with high probability are enforced to match the given text embedding. This
 123 gradient update method is named SDS, which is formulated as follows:

$$\nabla_\xi \mathcal{L}_{sds}(\xi) = \mathbb{E}_{t,\epsilon} \left[w(t) (\epsilon_\phi(\mathbf{x}_t; \mathbf{y}, t) - \epsilon) \frac{\partial \mathbf{x}_t}{\partial \xi} \right] \quad (6)$$

124 where the noise $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ and $w(t)$ is a weighting function. The neural rendering pipeline $\Gamma(\xi)$
 125 and the diffusion model ϕ as modular components of the framework, are amenable to selection. This
 126 offers a practical way for text-to-3D synthesis.

127 **Neural Texture for 3D-Consistent Rendering.** As compared with previous works [15, 10, 9]
 128 on optimizing a NeRF model, we directly optimize a neural texture over the UV space. Direct
 129 texture optimization is crucial to enforce 3D consistency and preserve texture details. In detail,
 130 first we exploit the pre-trained latent diffusion model as the optimization guidance, where it has an
 131 encoder V_e , a latent diffusion model's denoiser ϵ_ϕ , and a decoder V_d . In Decorate3D , the neural
 132 texture $\psi \in \mathbb{R}^{H \times W \times 3}$ is set to 3 feature channels to fit the requirement of the diffusion model.
 133 The gradient calculated via SDS is backpropagated through the diffusion model's encoder to the
 134 UV neural texture at a high resolution of 512×512 . Albeit yielding high-resolution images, the
 135 optimization computation is reasonable since the optimization exerts an effect on the latent code
 136 $\mathbf{z}_t^\psi = V_e(\mathcal{R}(\psi, \mathcal{M}, \mathcal{P}))$ with a resolution of 64×64 , where V_e is the diffusion model's encoder and
 137 \mathcal{P} is the sampled camera pose. All network parameters are fixed, while only the neural texture ψ is
 138 trainable.

139 Here, we call the optimized texture *UV Neural Texture* by following the deferred neural rendering
 140 [30]. The optimized texture does not act on RGB space, and they are neural features that need a
 141 *Neural Renderer* to interpret, which will be introduced shortly.

142 **Structure-Aware SDS Optimization.** Apart from addressing the 3D-consistency issues, we also
 143 want to maintain the coherence between 3D geometry and the generated texture. For example, for a
 144 3D human model, the generated facial texture should be attached to the 3D region corresponding

145 to the face. This structure-aware requirement is overlooked by the naive SDS optimization leading
 146 to a multi-face Janus problem. To moderate this issue, Decorate3D utilizes structure-aware SDS
 147 optimization by exploiting a depth-guided latent diffusion model [28, 22].

148 In detail, at each denoising step t , we compute the image latent code \mathbf{z}_t^ψ produced by the encoder V_e
 149 and the depth latent code $\mathbf{z}_t^{\psi_d}$ downsampled from the estimated depth map. They are concatenated
 150 as \mathbf{z}_t and forwarded to the latent diffusion module (Eq. 8). Then, following the SDS technique, the
 151 gradient of neural texture is computed as follows (Eq. 9):

$$\mathbf{z}_t = [\mathbf{z}_t^\psi, \mathbf{z}_t^{\psi_d}] \quad (7)$$

$$\tilde{\epsilon}_\phi(\mathbf{z}_t; \mathbf{y}, t) = \epsilon_\phi(\mathbf{z}_t; t) + \lambda[\epsilon_\phi(\mathbf{z}_t; \mathbf{y}, t) - \epsilon_\phi(\mathbf{z}_t; t)] \quad (8)$$

$$\nabla_\psi \mathcal{L}_{sds}(\phi, \psi) = \mathbb{E}_{t,\epsilon} \left[w(t)(\tilde{\epsilon}_\phi(\mathbf{z}_t; \mathbf{y}, t) - \epsilon_\phi(\mathbf{z}_t; t)) \frac{\partial \mathbf{z}_t}{\partial V_e} \frac{\partial V_e}{\partial \psi} \right] \quad (9)$$

152 where ϵ_ϕ is the diffusion model’s denoiser, \mathbf{y} is the text embedding from the Transformer, and $w(t)$
 153 is a weighting function. Decorate3D uses the classifier-free guidance scheme [14] as stated in Eq. 8,
 154 and the guidance weight λ for text conditioning is set to 100. To match the depth-guided diffusion
 155 model, depth maps are predicted from the rendered views by using the depth estimator [20] of the
 156 depth-guided diffusion model rather than the depth buffer of the rendering pipeline. The overall SDS
 157 optimization framework is shown at the left bottom of Fig. 2.

158 In addition to structure-aware SDS optimization, Decorate3D adopts structure-aware initialization.
 159 Thanks to the decomposition phase, Decorate3D has an inherent texture that well matches the
 160 geometry information extracted from the real-world object. Therefore, Decorate3D can initialize the
 161 neural texture ψ with the output of MLP_D . This initialization plays a vital role in the optimization
 162 process, which can accelerate the convergence speed as well as provide structure-aware information
 163 to help it converge to a structure-aware solution. We verify the effectiveness of the structure-aware
 164 techniques in Sec 4.4.

165 **Decorate3D’s Neural Renderer.** As aforementioned, directly optimizing the UV texture with
 166 SDS does not yield a traditional RGB UV texture for the rendering pipeline. Fig. 3 shows an
 167 example of SDS optimization for text-to-image generation. The optimized texture, *i.e.* neural
 168 texture, requires a neural interpreter to convert it back to RGB space after the SDS optimization.
 169 Previous work [10] applied a statistical linear transformation
 170 matrix to the optimized UV texture to interpret it in the RGB
 171 UV space. But this transformation matrix is suboptimal, and
 172 the misalignment will result in color shifts and artifacts.

173 In Decorate3D , we introduce a neural renderer as the neural
 174 interpreter to synthesize photo-realistic images from the optimized
 175 neural texture. This rendering process is formulated as
 176 follows:

$$I_{\mathcal{P}} = V(\mathcal{R}(\psi, \mathcal{M}, \mathcal{P})) \quad (10)$$

177 where ψ is the neural texture, V denotes the network of the
 178 neural renderer, and $I_{\mathcal{P}}$ is the rendered image given the camera
 179 pose \mathcal{P} . Here, Decorate3D ’s neural renderer network V is the encoder-decoder (*i.e.* VAE) module of
 180 the depth-guided latent diffusion model [22, 28]. The network V and the traditional renderer \mathcal{R} have
 181 amalgamated to form the neural renderer.

182 3.3 Few-View Resampling Training

183 With the neural texture, Decorate3D can render 3D-consistent high-quality views $I_{\mathcal{P}}$ through the
 184 Neural Renderer with Eq. 10 by sampling different camera poses. But there still exists a practical
 185 issue, *i.e.* jittering artifacts, that arises when rendering across different views. This is caused by the
 186 view-dependent neural renderer V . We empirically find that directly feeding the global UV neural
 187 texture ψ to the diffusion model’s VAE (*i.e.* $\mathcal{R}(V(\psi), \mathcal{M}, \mathcal{P})$) by swapping the network and renderer
 188 in Eq. 10) can produce the RGB UV texture for the traditional rendering pipeline. But this naive
 189 approach usually yields blurry texturing results, as demonstrated in Sec 4.4.

190 To solve the jittering problem, we devise a few-view resampling (FVR) training method, that transfers
 191 the synthesized view-dependent images into a global UV texture $\tilde{\psi}$. As shown in the right bottom of



Figure 3: The text prompt is “An apple tree”. The left shows the image from the direct SDS optimization on a 2D image. The right shows the image after applying interpreter.

192 Fig. 2, we use a $\text{MLP}_{\tilde{\psi}}$ to represent the UV texture $\tilde{\psi}$. In the FVR training, we sample N rendered
 193 views using the neural renderer with the neural texture ψ . The sampled N views should overlay
 194 the mesh surface as much as possible. But setting a big N may have negative effects, leading to
 195 over-smoothing textures on the overlapped areas that suffer jittering artifacts. The FVR training loss
 196 is defined as follows:

$$\mathcal{L}_{FVR}(\tilde{\psi}) = \frac{1}{N} \sum_i^N \left\| \mathcal{R}(\text{MLP}_{\tilde{\psi}}(\tilde{\mathbf{v}}), \mathcal{M}, \mathcal{P}_i) - V(\mathcal{R}(\psi, \mathcal{M}, \mathcal{P}_i)) \right\|^2 \quad (11)$$

197 where $\tilde{\mathbf{v}}$ denotes the positional encoding of the 2D UV coordinates of the RGB UV texture $\tilde{\psi}$, and
 198 \mathcal{P}_i denotes the i -th sampled camera pose. The pose is sampled in spherical coordinates, with two
 199 elevation angles θ_{cam} chosen from $\{-20^\circ, +20^\circ\}$, four azimuth angles β_{cam} uniformly sampled
 200 between $[0^\circ, 360^\circ]$, and an appropriate viewing distance r_{cam} . In our experiments, N is set to 8.

201 3.4 Super-Resolution on UV Texture

202 From the SDS approach, we can only synthesize images at a resolution of 512×512 . To further
 203 improve the spatial resolution of the UV texture $\tilde{\psi}$, Decorate3D applies a super-resolution (SR)
 204 diffusion model with a $\times 4$ scale factor [24, 29] on the UV texture to obtain a 2048×2048 UV texture.
 205 The rationale behind this successful application of UV texture upscaling is that the upsampling
 206 operation has a spatial locality, concentrating on local textures such as edges. For this reason, SR
 207 models are usually trained on cropped image patches [27, 4] instead of the whole image to increase
 208 the training efficiency. This spatial locality of SR allows the SR model trained on natural images to
 209 be directly applied to the UV texture. Directly applying the SR operation to the UV texture is free
 210 from any jittering issues, and 3D consistency is well preserved.



211 Figure 4: Text-driven texture generation results of Decorate3D .

211 4 Experiments

212 We compare Decorate3D to the state-of-the-art (SOTA) techniques for text-to-3D texture generation
 213 and evaluate its performance from both qualitative and quantitative perspectives. The 3D assets and
 214 extended videos are presented in the supplementary material.

215 **4.1 Implementation Details**

216 To evaluate our approach, we collect real-world datasets from 14 different objects that vary in
 217 complexity, including boxes, monitors, shoes, statues, dolls, humans, and lighthouses. Ten objects
 218 are captured in the wild using a smartphone, and four objects are selected from some public real-
 219 world datasets [33, 8]. The size of each dataset ranges from 70-300 images, and the images are
 220 downsampled to a resolution width of 640. We use the Adam optimizer to optimize the $\text{MLP}_{\mathcal{D}}$,
 221 MLP_S and $\text{MLP}_{\tilde{\psi}}$ with a learning rate of 1×10^{-3} , and the ψ with a learning rate of 1×10^{-2} . The
 222 neural texture optimization in the decoration phase takes about 2 hours for 100K iterations, and the
 223 FVR training takes about 5 minutes for 30K iterations, which are measured on 8 NVIDIA V100
 224 GPUs with batch size 1. Please refer to the supplementary material for more details.

225 **4.2 Qualitative Evaluation and Comparison**



Figure 5: Qualitative comparison with other text-driven texture generation methods, including CLIP-Mesh [7], Latent-Paint [10] and TEXTure [21]. Results from DreamFusion [15] are also shown here. Please zoom in for better visualization.

226 **3D-Consistent Text-Driven Texture Generation.** In Fig. 4, we show the textured mesh and its
 227 corresponding real-world image. It can be observed that Decorate3D is able to produce high-quality
 228 textured mesh. As the mesh geometry is fixed, we find only texture-related keywords such as nouns
 229 or adjectives will affect the generated results, highlighted with a purple color in the prompt. For
 230 instance, in the prompt “An astronaut stands up in the milky way”, “astronaut” and “milky way”
 231 dominate the generation, but the verb “stands up” and the quantifier “an” do not affect the results.

232 Fig. 5 shows the qualitative comparison with SOTA text-driven texture generation methods, including
 233 CLIP-Mesh [7], Latent-Paint [10] and TEXTure [21]. They all took the reconstructed 3D mesh as
 234 input, and the geometry was fixed during texture generation. We can observe that the generated
 235 textures from CLIP-Mesh and Latent-Paint lack texture details. The concurrent work TEXTure has
 236 some competitive results compared to Decorate3D, such as the results for shoes, because they use
 237 a depth-guidance strategy similar to ours. But TEXTure produces weird artifacts on the complex
 238 surfaces. For example, in the case of humans, the generated textures from TEXTure have obvious
 239 seams, and textures from different parts seem to be glued together and look messy. The artifacts are
 240 caused by the iterative update strategy across views. By contrast, the results of Decorate3D have
 241 fruitful and clear details. We also present the results from DreamFusion [15]. It tends to produce
 242 blurry results that do not actually match the text prompt.

243 **3D-Consistent Texture Editing.** Thanks to the decomposed 3D representations, Decorate3D opens
 244 up an easy way to edit textures. Fig. 6 depicts a practical and distinctive use case that can be

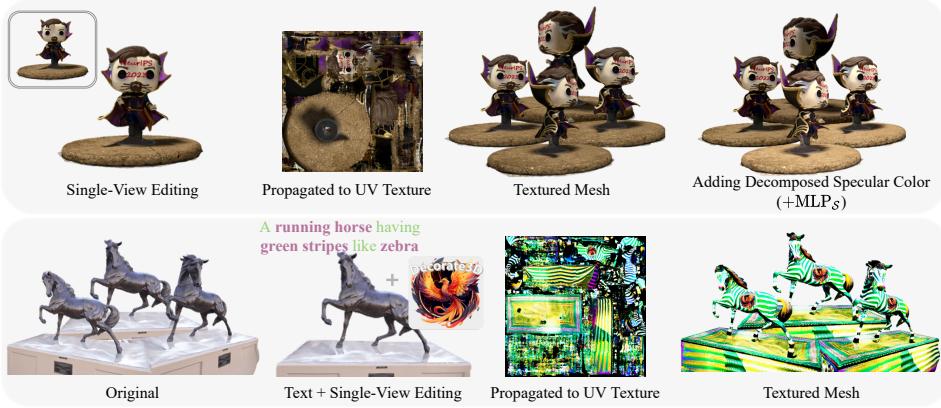


Figure 6: Two examples of applying Decorate3D for 3D-consistent texture editing. As showcased in the first example, a scribble of ‘NeurIPS 2023’ pattern is painted onto a 2D image (aka. a single rendered view), which is propagated by Decorate3D to the UV texture, updating the 3D-consistent rendered output.

Table 1: Evaluating the correlation of text-driven generated results with their text input using different CLIP models. The CLIP-Mesh’s scores may overfit, as it uses the same CLIP model for training and eval. The NIQE is a no-reference image quality evaluation metric.

Method	R-Precision ↑			
	CLIP B/32	CLIP B/16	CLIP L/14	NIQE ↓
CLIP-Mesh [7]	31.95 ± 3.01	30.16 ± 2.63	23.99 ± 3.71	16.28 ± 0.94
DreamFusion [15]	29.65 ± 4.71	29.63 ± 5.11	24.38 ± 4.46	16.68 ± 0.47
Latent-Paint [10]	25.18 ± 4.30	26.19 ± 2.81	21.22 ± 3.03	17.71 ± 1.39
TEXTure [21]	29.68 ± 3.56	28.65 ± 2.71	23.15 ± 3.31	14.83 ± 0.78
Decorate3D (Ours)	30.42 ± 3.47	30.47 ± 3.07	24.89 ± 3.41	14.82 ± 0.66

245 accomplished by Decorate3D . We edit one of the rendered images (not the UV texture), and
246 propagate the editing to the UV texture to synthesize 3D-consistent edited textures from different
247 views. The decomposed view-dependent specular texture can be added back to the textured model.
248 Additionally, Decorate3D allows secondary editing of the text-driven texture generation.

249 4.3 Quantitative Evaluation and Comparison

250 **CLIP R-Precision Metric.** Following DreamFusion, we evaluate the CLIP R-Precision [17], an
251 automated metric for the consistency of rendered multi-view images with respect to the input prompt.
252 Here, we calculate the average CLIP score from the front-side, left-side, right-side, and back-side
253 views. We use 70 prompts of 14 objects to generate the test results. For visual quality, we measure
254 the NIQE [12] on rendered images, which is a no-reference image quality assessment. Tab. 1 reports
255 the CLIP R-Precision scores of the compared methods.

256 **User Studies.** We invite 44 volunteers to evaluate Decorate3D and
257 its competitors using the mean-opinion-score (MOS) test. For each
258 question, we prepare 5 video results, and the participants are asked to
259 rate the results on a scale from 1 (worst) to 5 (best) based on the overall
260 visual quality of the results and the degree to which they match the text
261 prompt. In the end, we receive 1100 responses from the 44 volunteers.
262 Fig. 7 shows the average MOS scores of the compared methods. As
263 can be seen, Decorate3D is shown to be more favored by human users.

264 4.4 Ablation Studies

265 **Initialization of UV Neural Texture.** Fig. 9a shows an ablation study
266 on the initialization of the UV neural texture. As observed, the optimization initialized with MLP_D
267 converges much faster than the random initialization. Moreover, the MLP_D provides a very strong
268 optimization prior to helping the SDS optimization converge to a better solution that well matches
269 the mesh geometry. For example, initialization using MLP_D can generate the correct facial texture
270 that fits the geometry, but on the contrary, the random initialization fails.

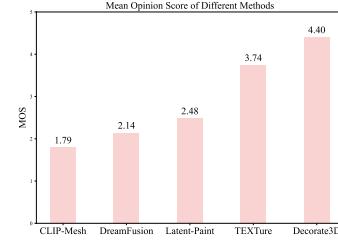


Figure 7: User study results gathered from 44 participants. The results are the average of all responses.

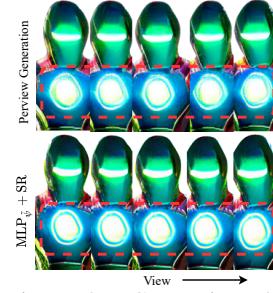


Figure 8: Comparison between the per-view generation and rendering with MLP_ψ .

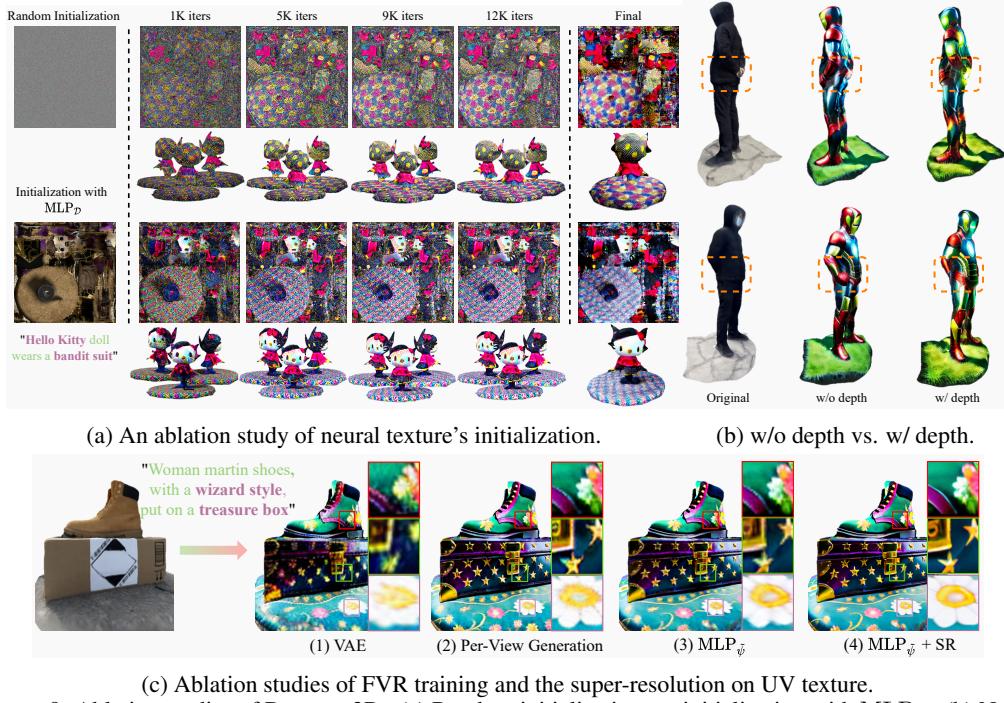


Figure 9: Ablation studies of Decorate3D . (a) Random initialization vs. initialization with $\text{MLP}_{\mathcal{D}}$. (b) Naive SDS without depth vs. structure-aware SDS guided by depth. As observed, the naive SDS cannot generate the geometry-matching texture for the man’s hands. (c) We compare the rendered results using (1) UV texture yielded by the diffusion model’s VAE, (2) Per-view generation, (3) FVR trained $\text{MLP}_{\tilde{\psi}}$, and (4) $\text{MLP}_{\tilde{\psi}} + \text{SR}$.

271 **Effectiveness of Structure-Aware SDS.** Fig. 9b presents the difference between the naive SDS
 272 without depth guidance and the structure-aware SDS with depth guidance. We can observe that the
 273 structure-aware SDS is able to produce geometry-matching textures. For example, the human’s arms
 274 should be clasped behind the back rather than being akimbo.

275 **Is the Few-View Resampling Training Necessary?** Fig. 8 compares
 276 the per-view generated results of $V(\mathcal{R}(\psi, \mathcal{M}, \mathcal{P}))$ with the results rendered
 277 using the FVR-trained UV texture. As observed, results of per-view
 278 generation have jittering effects between different views (look at the circle
 279 pattern), but results with $\text{MLP}_{\tilde{\psi}}$ have achieved 3D consistency. Fig. 9c (1)
 280 presents the results using an alternative solution to eliminate the jittering,
 281 i.e. directly feeding the neural texture into the diffusion model’s VAE
 282 by $\mathcal{R}(V(\psi), \mathcal{M}, \mathcal{P})$. As observed, following this way yields blurry tex-
 283 tures. Given an RGB UV texture $\tilde{\psi}$, Decorate3D can directly apply SR
 284 on $\tilde{\psi}$. By comparing Fig. 9c (3) and Fig. 9c (4) we can observe more
 285 clear results after SR. Fig. 10 shows the effectiveness of FVR training’s
 286 hyperparameter N . Setting a big N for FVR will cause blurry results.

287 5 Conclusion

288 Decorate3D offers a practical way of decorating real-world 3D models with a user-friendly approach
 289 through text-driven texture generation. Our Decorate3D ’s techniques, combining directly optimizing
 290 UV neural texture, structure-aware optimization, FVR training, and SR enhancement on UV texture,
 291 collectively advance the line of 3D texture generation in pursuit of the best possible quality. Extensive
 292 experiments demonstrate the superiority of Decorate3D over SOTA methods.

293 **Limitations** First, the style and quality of the generated texture heavily depend on the pre-trained
 294 stable diffusion models. Second, even though Decorate3D adopts structure-aware optimization tech-
 295 niques, the multi-face Janus problem still remains in a flat geometry. For example, Decorate3D cannot
 296 distinguish the front and back sides of a monitor object, as shown in our supplementary material.
 297 Lastly, Decorate3D does not support jointly optimizing both the geometry and texture, which may
 298 result in inconsistency with the mesh geometry. We remain these problems for future study.

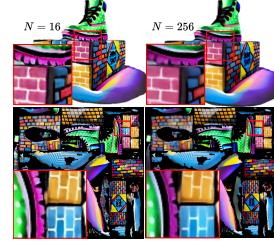


Figure 10: Results with different N for FVR. A bigger N leads to blurry textures.

299 **References**

- 300 [1] Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow
301 image editing instructions. In *The Conference on Computer Vision and Pattern Recognition*
302 (*CVPR*), 2023.
- 303 [2] Dave Zhenyu Chen, Yawar Siddiqui, Hsin-Ying Lee, Sergey Tulyakov, and Matthias Nießner.
304 Text2tex: Text-driven texture synthesis via diffusion models, 2023.
- 305 [3] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and
306 appearance for high-quality text-to-3d content creation. *arXiv preprint arXiv:2303.13873*, 2023.
- 307 [4] Yanhui Guo, Xiaolin Wu, and Xiao Shu. Data acquisition and preparation for dual-reference
308 deep learning of image super-resolution. *IEEE Transactions on Image Processing*, 31:4393–
309 4404, 2022.
- 310 [5] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim
311 Salimans. Cascaded diffusion models for high fidelity image generation. *Journal of Machine*
312 *Learning Research*, 23(47):1–33, 2022.
- 313 [6] Heewoo Jun and Alex Nichol. Shap-e: Generating conditional 3d implicit functions. *arXiv*
314 *preprint arXiv:2305.02463*, 2023.
- 315 [7] Nasir Mohammad Khalid, Tianhao Xie, Eugene Belilovsky, and Tiberiu Popa. CLIP-mesh:
316 Generating textured meshes from text using pretrained image-text models. In *SIGGRAPH Asia*
317 *2022 Conference Papers*. ACM, nov 2022.
- 318 [8] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Bench-
319 marking large-scale scene reconstruction. *ACM Transactions on Graphics*, 36(4), 2017.
- 320 [9] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten
321 Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d
322 content creation. In *IEEE Conference on Computer Vision and Pattern Recognition* (*CVPR*),
323 2023.
- 324 [10] Gal Metzger, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. Latent-nerf for
325 shape-guided generation of 3d shapes and textures. *IEEE Conference on Computer Vision and*
326 *Pattern Recognition* (*CVPR*), 2023.
- 327 [11] Ben Mildenhall, Pratul Srinivasan, Matthew Tancik, Jonathan Barron, Ravi Ramamoorthi, and
328 Ren Ng. Nerf: representing scenes as neural radiance fields for view synthesis. In *European*
329 *Conference on Computer Vision (ECCV)*, 2020.
- 330 [12] Anish Mittal, Rajiv Soundararajan, and Alan C. Bovik. Making a “completely blind” image
331 quality analyzer. *IEEE Signal Processing Letters*, 20(3):209–212, 2013.
- 332 [13] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics
333 primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, July
334 2022.
- 335 [14] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin,
336 Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image genera-
337 tion and editing with text-guided diffusion models. In *International Conference on Machine*
338 *Learning*, pages 16784–16804. PMLR, 2022.
- 339 [15] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using
340 2d diffusion. *International Conference on Learning Representations (ICLR)*, 2023.
- 341 [16] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
342 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
343 models from natural language supervision. In *International conference on machine learning*,
344 pages 8748–8763. PMLR, 2021.
- 345 [17] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-
346 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya
347 Sutskever. Learning transferable visual models from natural language supervision. In Marina
348 Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine*
349 *Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR,
350 18–24 Jul 2021.

- 351 [18] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical
 352 text-conditional image generation with clip latents, 2022.
- 353 [19] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark
 354 Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on*
 355 *Machine Learning*, pages 8821–8831. PMLR, 2021.
- 356 [20] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards
 357 robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE*
 358 *transactions on pattern analysis and machine intelligence*, 44(3):1623–1637, 2020.
- 359 [21] Elad Richardson, Gal Metzer, Yuval Alaluf, Raja Giryes, and Daniel Cohen-Or. Texture:
 360 Text-guided texturing of 3d shapes. *arXiv preprint arXiv:2302.01721*, 2023.
- 361 [22] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
 362 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF*
 363 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, June
 364 2022.
- 365 [23] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed
 366 Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim
 367 Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image
 368 diffusion models with deep language understanding, 2022.
- 369 [24] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad
 370 Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis*
 371 *and Machine Intelligence*, 2022.
- 372 [25] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In
 373 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- 374 [26] Junyoung Seo, Wooseok Jang, Min-Seop Kwak, Jaehoon Ko, Hyeonsu Kim, Junho Kim, Jin-
 375 Hwa Kim, Jiyoung Lee, and Seungryong Kim. Let 2d diffusion model know 3d-consistency for
 376 robust text-to-3d generation. *arXiv preprint arXiv:2303.07937*, 2023.
- 377 [27] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop,
 378 Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an
 379 efficient sub-pixel convolutional neural network. In *Proceedings of the IEEE conference on*
 380 *computer vision and pattern recognition*, pages 1874–1883, 2016.
- 381 [28] Stability-AI. Depth-guided stable diffusion model. <https://github.com/Stability-AI/stablediffusion>, 2023.
- 383 [29] Stability-AI. Diffusion models for super-resolution. <https://huggingface.co/CompVis/ldm-super-resolution-4x-openimages>, 2023.
- 385 [30] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image
 386 synthesis using neural textures. *Acm Transactions on Graphics (TOG)*, 38(4):1–12, 2019.
- 387 [31] Can Wang, Ruixiang Jiang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao.
 388 Nerf-art: Text-driven neural radiance fields stylization. *arXiv preprint arXiv:2212.08070*, 2022.
- 389 [32] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang.
 390 Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction.
 391 *Advances in Neural Information Processing Systems*, 34:27171–27183, 2021.
- 392 [33] Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan Ren, Lei Zhou, Tian Fang, and Long
 393 Quan. Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. In
 394 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,
 395 June 2020.
- 396 [34] Jonathan Young. Xatlas: Mesh parameterization / uv unwrapping library. <https://github.com/jpcy/xatlas>, 2023.
- 398 [35] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion
 399 models, 2023.