

# SceneTex: High-Quality Texture Synthesis for Indoor Scenes via Diffusion Priors

Dave Zhenyu Chen<sup>1</sup> Haoxuan Li<sup>1</sup> Hsin-Ying Lee<sup>2</sup> Sergey Tulyakov<sup>2</sup> Matthias Nießner<sup>1</sup>  
<sup>1</sup>Technical University of Munich <sup>2</sup>Snap Research  
<https://daveredrum.github.io/SceneTex/>



Figure 1. We introduce SceneTex, a text-driven texture synthesis architecture for 3D indoor scenes. Given scene geometries and text prompts as input, SceneTex generates high-quality and style-consistent textures via depth-to-image diffusion priors.

## Abstract

We propose SceneTex, a novel method for effectively generating high-quality and style-consistent textures for indoor scenes using depth-to-image diffusion priors. Unlike previous methods that either iteratively warp 2D views onto a mesh surface or distillate diffusion latent features without accurate geometric and style cues, SceneTex formulates the texture synthesis task as an optimization problem in the RGB space where style and geometry consistency are properly reflected. At its core, SceneTex proposes a multiresolution texture field to implicitly encode the mesh appearance. We optimize the target texture via a score-distillation-based objective function in respective RGB renderings. To further secure the style consistency across views, we introduce a cross-attention decoder to predict the RGB values by cross-attending to the pre-sampled reference locations in each instance. SceneTex enables various and accurate texture synthesis for 3D-FRONT scenes, demonstrating significant improvements in visual quality and prompt fidelity over the prior texture generation methods.

## 1. Introduction

Synthesizing high-quality 3D contents is an essential yet highly demanding task for numerous applications, including gaming, film making, robotic simulation, autonomous driving, and upcoming VR/AR scenarios. With an increasing number of 3D content datasets, the computer vision and graphics community has witnessed a soaring research interest in the field of 3D geometry generation [2, 12, 36, 38, 40, 60, 68, 73]. Despite achieving a remarkable success in 3D geometry modeling, generating the object appearance, i.e. textures, is still bottlenecked by laborious human efforts. It typically requires a substantially long time for designing and adjustment, and immense 3D modelling expertise with tools such as Blender. As such, automatic designing and augmenting the textures has not yet been fully industrialized due to a huge demand for human expertise and financial expenses.

Leveraging the recent advances of 2D diffusion models, tremendous progress has been made for text-to-3D generation, especially for synthesizing textures of given shapes [8, 39, 50]. Seminal work such as Text2Tex [8] and Latent-Paint [39] have achieved great success in generating high-quality appearances for objects, facilitating high-fidelity texture synthesis from input prompts. Despite the

fascinating results on objects, upscaling these methods to generating textures for an entire scene still confronts several challenges. On one hand, methods that autoregressively project 2D views to 3D object surface [8, 50] usually suffer from texture seams, accumulated artifacts, and loop closure issues. It is also quite difficult to maintain style consistency in the scene if every object is textured individually. On the other hand, score-distillation-based approaches [39] perform texture optimization in the low-resolution latent space, often resulting in blurry RGB textures and incorrect geometry details. As such, previous text-driven attempts fail to deliver high-quality textures for 3D scenes.

To address the aforementioned challenges, we propose SceneTex, a novel architecture to generate high-quality and style-consistent texture for indoor scene meshes by leveraging depth-to-image diffusion priors. Unlike previous methodologies that iteratively warp 2D views onto mesh surfaces, we take a different approach by framing the texture synthesis as a texture optimization task in RGB space via diffusion priors. At its core, we introduce a multiresolution texture field to implicitly represent the appearance of the mesh. To faithfully represent the texture details, we adopt a multiresolution texture to store texture features at multiple scales. This enables our architecture to flexibly learn both low and high frequency appearance information. To secure the style consistency of the generated texture, we incorporate a cross-attention decoder to reduce style incoherence introduced by self-occlusion. Concretely, each decoded RGB values are produced by cross-attending to the pre-sampled reference surface locations scattered across each object. This way, we further secure the global style consistency within each instance, as every visible location receives a global reference to the whole instance appearance.

We show that SceneTex has the capacity to facilitate versatile and accurate texture synthesis for indoor scenes with given language cues. We demonstrate in extensive experiments that SceneTex places a strong emphasis on both style and geometry consistency. The proposed method performs favorably against other text-driven texture synthesis methods in terms of 2D metrics such as CLIP score [48] and Inception Score [60], and user study on a subset of the 3D-FRONT dataset [20].

We summarize our technical contributions as follows:

- We design a novel framework for generating high-quality scene textures in high resolution using depth-to-image diffusion priors.
- We propose an implicit texture field to encode the object appearance at multiple scales, leveraging a multiresolution texture to faithfully represent rich texture details.
- We incorporate a cross-attention texture decoder to secure the global style-consistency for each instance, producing more visually appealing and style-consistent textures for

3D-FRONT scenes compared against previous synthesis methods.

## 2. Related work

**Feed-forward 3D Generation.** The advancement of 3D generation has adhered to the progress of 2D generative models. Adopting different backbone techniques, from variational autoencoders [34], generative adversarial networks [24], autoregressive transformers [64], to the recent diffusion models [19, 27], 3D models have been trained on 3D data of various representations, including voxels [10, 36, 55, 60, 68], point clouds [2, 38, 74], meshes [46, 73], signed distance functions [12, 14, 15, 17, 40], and more. However, unlike the ubiquity of 2D images and videos, 3D data is inherently scarce and poses significant challenges in terms of acquisition and annotation. Recent efforts have sought to address this issue by utilizing differentiable rendering techniques to learn from 2D images [1, 3, 4, 21, 25, 41, 44, 53, 54, 56, 59, 69, 70]. Although these models typically demonstrate proficiency in specific shape categories, they are incapable of handling 3D generation from free-form texts.

**3D Generation with 2D diffusion models.** Recently, significant strides have been made in the field of vision-language integration [5–7, 13, 31, 48, 57, 66]. The advancements in text-to-image generation, particularly diffusion models trained on large-scale image collections [19, 27, 28, 43, 52], have prompted the integration of pretrained 2D diffusion models as priors to facilitate 3D generation. Two main streams of work have emerged. The first branch directly incorporates the output of the 2D diffusion models along with the depth information. TEXTure [50] and Text2Tex [8] perform texturing on given meshes with a depth-aware variation of diffusion models. Other methods generate 3D scenes, where the geometry information is either jointly predicted [61] or obtained from off-the-shelf depth estimator [30]. The second branch of methods [9, 35, 39, 47, 65, 67] attempt to distill knowledge from pretrained 2D diffusion models with the Score Distillation Sampling (SDS) [47] technique and its subsequent improved variations in a per-prompt optimization manner. In contrast, we take advantages of the distilled knowledge from the depth-conditioned 2D diffusion priors to enable high-quality 3D texture synthesis.

**3D Scene Texturing.** In this work, we focus on generating high-quality textures for 3D scenes. 3D scene texturing has been studied by applying the 2D style transfer techniques [22, 23, 33] to 3D domain [11, 16, 26, 29, 32, 71]. However, these methods often emphasize low-level styles without semantic understanding. While existing 3D generation methods leveraging 2D diffusion models can theoretically be applied to 3D scene texturing, those based on inpainting [8, 50] suffer from visible seams and accumulated artifacts, while distillation-based methods [39] often

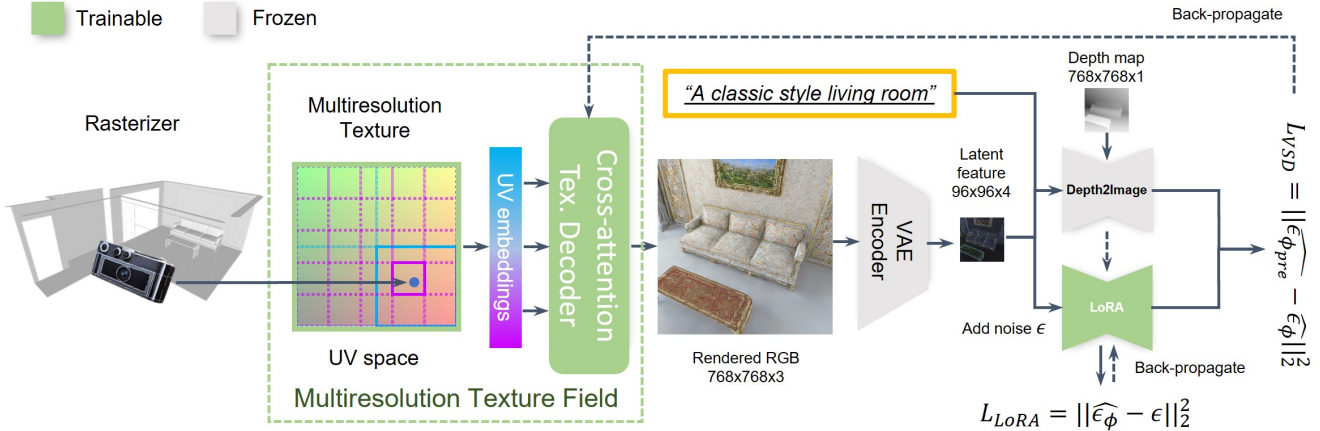


Figure 2. **Texture synthesis pipeline.** The target mesh is first projected to a given viewpoint via a rasterizer [37]. Then, we render an RGB image with the proposed multiresolution texture field module. Specifically, each rasterized UV coordinate is taken as input to sample the UV embeddings from a multiresolution texture. Afterward, the UV embeddings are mapped to an RGB image of shape  $768 \times 768 \times 3$  via a cross-attention texture decoder. We use a pre-trained VAE encoder to compress the input RGB image to a  $96 \times 96 \times 4$  latent feature. Finally, the Variational Score Distillation loss [67] is computed from the latent feature to update the texture field.

produce blurry textures with incorrect geometry details. In contrast, we optimize the target scene texture with accurate geometric cues and decode the high-resolution scene appearance via the proposed multiresolution texture field module, facilitating 3D scene texture synthesis with much better visual quality.

### 3. Method

The objective of our work is to texture an entire 3D scene with diffusion priors as the critic. In this section, we begin by introducing a Multiresolution Texture Field module to produce high-quality RGB textures, which consists of two key components: Multiresolution Texture and Cross-attention Texture Decoder. The Multiresolution Texture is integrated to faithfully represent both the low- and high-frequency texture details at various scales (Sec. 3.1). Subsequently, to tackle the style-inconsistency issue brought by limited field of view and self-occlusion, the Cross-attention Texture Decoder module is incorporated to enforce a global style-awareness for each object in the scene (Sec. 3.2). Finally, we adopt a pretrained diffusion model as training critic to dynamically distillate realistic scene appearance from the 2D depth-conditioned diffusion priors. (Sec. 3.3). The entire synthesis architecture is presented in Fig. 2.

#### 3.1. Multiresolution Texture Field

The core of texture synthesis with 2D priors lies in generating RGB values visible to a series of pre-defined viewpoints. Previous methods maintain a  $64 \times 64 \times 4$  latent map and operate directly on it with the SDS loss [39, 47]. This latent map is decoded via the variational autoencoder of the pre-trained diffusion model after convergence. The

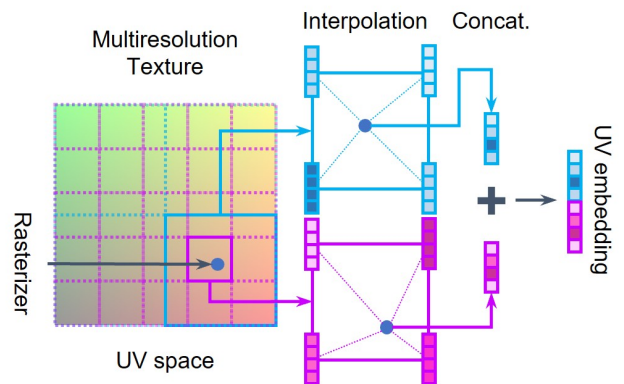


Figure 3. **Multiresolution Texture.** We use a multiresolution feature grid to encode positional features at different scale in the UV space. For a query UV coordinate, we interpolate the grid features at respective resolutions. The interpolated grid features are concatenated as the final UV embedding for the query UV coordinate.

optimization process is technically view-consistent, as the diffusion priors are leveraged from numerous perspectives. Notwithstanding, we observe that the decoded RGB textures often carry patch-like artifacts and are subsequently inconsistent with the given geometry. This is caused by the mismatches between the low-resolution latent map and high-resolution RGB images, and the lack of perspective transformation of the same latent code in different views.

To tackle those inherent disadvantages of representing the target texture via a low-resolution latent map, we adopt an implicit texture field that queries the texture features with given UV coordinates. At its core, we integrate a multiresolution texture to prevent oversimplified appearance with-



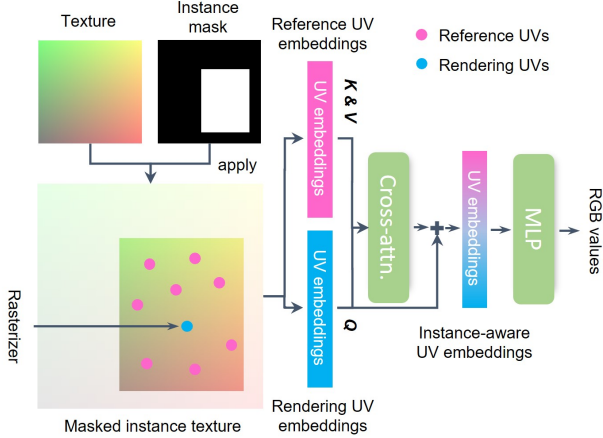


Figure 4. **Cross-attention Texture Decoder.** For each rasterized UV coordinate, we apply a UV instance mask to mask out the corresponding instance texture features. Then, we obtain the rendering UV embeddings for the rasterized locations in the view. At the same time, we extract the texture features for the pre-sampled UVs scattered across this instance as the reference UV embeddings. We deploy a multi-head cross-attention module to produce the instance-aware UV embeddings. Here, we treat the rendering UV embeddings as the Query, and the reference UV embeddings as the Key and Value. Finally, a shared MLP maps the instance-aware UV embeddings to RGB values in the rendered view.

out any texture details. In particular, as shown in Fig. 3, we encode texture features for all query locations  $q$  at each scale, and concatenate those features as the output UV embeddings  $\mathcal{E}(q)$  to faithfully represent all texture details. The UV embeddings are then decoded to the final RGB texture by the cross-attention texture decoder introduced in the next section.

### 3.2. Cross-attention Texture Decoder

Since the texture is optimized in image space, instance textures are often constrained by limited field of view and self-occlusion. As a result, the optimized texture often suffers from style-inconsistency. Therefore, we propose a simple yet effective rendering module with global instance awareness to predict RGB values from UV embeddings. This is done by incorporating a multi-head cross-attention module to the texture features. As Fig 4 illustrates, for each rasterized UV coordinate, we apply a UV instance mask to mask out the corresponding instance texture features. Then, we obtain the rendering UV embeddings for the rasterized locations in the view. At the same time, we extract the texture features for the pre-sampled UVs scattered across this instance as the reference UV embeddings. We deploy a multi-head cross-attention module to produce the instance-aware UV embeddings. Here, we treat the rendering UV embeddings as the Query, and the reference UV embeddings as

the Key and Value. Finally, a shared MLP maps the global-aware UV embeddings to RGB values in the rendered view. We denote the whole rendering process as  $\mathcal{C} = f(\mathcal{E}(q); \theta)$ , where  $\mathcal{C}$  represents an RGB image at arbitrary resolution,  $f(\theta)$  is a differentiable function resembles the entire texture field with trainable parameters  $\theta$ .

### 3.3. Texture Field Optimization via VSD

We adopt a pre-trained ControlNet model as a critic to optimize the texturing module  $f(\theta)$  following the strategy of Latent-Paint [39], as shown in Fig. 2. Here, the UNet of a pre-trained latent diffusion model (LDM) [51] applied to calculate the gradients based on a low-resolution  $96 \times 96$  latent map. We observe that such low-resolution rendering often lead to broken visual quality and unsatisfactory view consistency. This is primarily due to the size mismatches between the  $96 \times 96$  optimization target and the final  $768 \times 768$  RGB output. Additionally, prior work exclude geometric cues from the diffusion priors, resulting in poor consistency between the generated textures and target geometry. To address those issues, we directly render an  $768 \times 768$  RGB image via querying the texture field  $\mathcal{C} = f(\mathcal{E}(q); \theta)$ . In each iteration, we first optimize  $\mathcal{C}$  via the VSD objective [67] with a pre-trained frozen depth-conditioned diffusion prior  $\phi_{\text{pre}}$  and a trainable LoRA module  $\phi$ :

$$\nabla_{\theta} \mathcal{L}_{\text{VSD}}(\theta) \approx \mathbb{E}_{t, \epsilon} [w(t) (\epsilon_{\phi_{\text{pre}}} - \epsilon_{\phi}) \frac{\partial f(\theta)}{\partial \theta}] \quad (1)$$

where  $\epsilon_{\phi_{\text{pre}}} = \phi_{\text{pre}}(f(\theta); y, d, t)$  and  $\epsilon_{\phi} = \phi(f(\theta); y, d, t)$ . We draw time step randomly by  $t \sim \mathcal{U}(0.02, 0.98)$ . The injected noise is  $\epsilon \sim \mathcal{N}(0, 1)$ .  $d$  is the depth map in the current viewpoint produced by the rasterizer.  $y$  is the noised input to the UNet. The weighting function  $w(t)$  is empirically set as  $w(t) = \sqrt{1 - \prod_{s=1}^t \alpha_s}$ . Note that  $\phi_{\text{pre}}$  and  $\phi$  are kept frozen when updating the parameters of the texture field  $\theta$ . After  $\mathcal{C} = f(\mathcal{E}(q); \theta)$  is updated via the gradients of VSD, we unfreeze and update the LoRA module  $\phi$ :

$$\mathcal{L}_{\text{LoRA}}(\phi) = \min_{\phi} \sum_{i=1}^n \mathbb{E}_{t, \epsilon} [\|\epsilon_{\phi}(f(\theta); y, d, t) - \epsilon\|_2^2] \quad (2)$$

### 3.4. Inference

Since the texture field  $f(\theta)$  only receives the UV coordinates as input, producing the final RGB texture is straightforward. For each position  $q_i$  in UV space, the corresponding pixel value  $c_i$  of the output RGB texture  $\mathcal{C}$  can be simply queried by  $c_i = f(\mathcal{E}(q_i); \theta)$ . Thanks to the multiresolution grid encoding, it is worth mentioning that there is no specification for the size of the final RGB texture, i.e. the resolution of the texture can be adjusted according to the computational resources.

## 4. Results

### 4.1. Implementation Details

We apply the ControlNet Depth model [72] for VSD optimization. The multiresolution texture is implemented by multiresolution hash encoding [42]. During each optimization iteration, we randomly pick 1 viewpoint from all perspective points scattered across the scene. We set the learning as 0.001 for optimizing the texture field and 0.0001 for fine-tuning the LoRA module. The entire optimization uses 5,000 viewpoints and takes 30,000 iterations to converge. To generate a more visually appealing appearance via VSD, We adopt the time annealing scheme following ProlificDreamer [67], where we sample time steps  $t \sim \mathcal{U}(0.02, 0.98)$  for the first 5,000 steps and then anneal into  $t \sim \mathcal{U}(0.02, 0.50)$  for the rest of the optimization. For the proposed cross-attention decoder, we pre-sampled 4,096 UV coordinates scattered across each instance. To enable cross-attention in such a long context, we implement the cross-attention module with Flash Attention v2 [18]. Each synthesis process takes around 20 hours to converge on an NVIDIA RTX A6000. After convergence, we generate a high-resolution 4,096 × 4,096 RGB image as the final scene texture. Our implementation uses the PyTorch [45] framework, with PyTorch3D [49] for rendering and texture projection.

### 4.2. Quantitative Analysis

We compare our method against texture synthesis methods appeared recently in Tab. 1, including Latent-Paint [39], MVDiffusion [63], and Text2Tex [8]. We experiment all methods on 10 3D-FRONT [20] scenes with 2 different text prompts for each scene. Here, we calculate CLIP score (CLIP) [48] and Inception Score (IS) [60] to measure the fidelity with input prompts and texture quality, respectively. Our method outperforms all baselines on the 2D automated metrics by a significant margin. We additionally report the User Study results from 75 participants about the Visual Quality (VQ) and Prompt Fidelity (PF) on a scale of 1-5. Our method is shown to be more favored by human users.

### 4.3. Qualitative Results

We show the qualitative comparisons in Fig. 5. Latent-Paint suffers from the over-saturation issue and hallucinates non-existing objects, such as the huge frame on the wall (see the first example in the first row in Fig. 5). Those unrealistic texture components are produced by the inaccurate geometric cues and the mismatch between the optimized latent representation and final texture. MVDiffusion [63] produces overall smooth but blurry and dimmed texture. It also fails to reflect the iconic properties in the prompts, such as “baroque” and “luxury”. Text2Tex [8] generates plausible textures for individual objects, but fails to achieve global

Method	2D Metrics		User Study	
	CLIP ↑	IS ↑	VQ ↑	PF ↑
Latent-Paint [39]	18.37	1.96	1.57	2.11
MVDiffusion [63]	18.47	2.83	3.09	3.12
Text2Tex [8]	20.83	2.87	2.62	3.04
SceneTex (Ours) w/o texture field	15.77	1.56	1.23	1.11
SceneTex (Ours) w/o multires. tex.	19.87	2.79	2.11	2.39
SceneTex (Ours) w/ cross-attn.	20.94	3.29	3.94	4.05
SceneTex (Ours)	<b>22.18</b>	<b>3.33</b>	<b>4.40</b>	<b>4.29</b>

Table 1. **Quantitative comparisons.** We report the 2D metrics and User Study results for quantitative comparisons, including: CLIP score (CLIP) [48], Inception Score (IS) [60], Visual Quality (Visual Quality), and Prompt Fidelity (PF). We show that our method produces textures with the highest quality.

style consistency across objects. In contrast, our method synthesizes high-quality with overall coherent styles within and across objects, reflecting the representative traits in the prompts with high-fidelity (see the baroque paintings above and the golden pillows below). We additionally visualize our texture synthesis results for different 3D-FRONT [20] scenes and input prompts in top-down views and close-up renderings in Fig. 6, further demonstrating the supreme texture quality and fidelity produced by our method.

### 4.4. Ablation Studies

We conduct ablation experiments on the key components of our method, including multiresolution texture (Sec. 3.1), and cross-attention decoder (Sec. 3.2). All comparisons are shown in Tab. 1 and Fig. 7.

**Does texture field produce better textures than RGB tensors?** Since only a few UVs are sampled by the rasterizer during each iteration, directly optimizing an RGB tensor as the output texture leads to noisy artifacts, as shown in the first column in Fig. 7. Additionally, the optimization is often difficult to converge with different gradient scales across the whole RGB tensor. As a result, the optimized texture appears to be broken and unrealistic. In contrast, the MLP in the proposed texture field module effectively smoothens the back-propagated gradients, producing much smoother and detailed textures.

**Does multiresolution texture improve the visual quality?** As previous studies indicate, implicit representations via MLPs tend to learn low-frequency information [58, 62]. As such, the texture field with single resolution produces an over-simplified appearance without texture details. Such texture lacks the characteristic properties of the input prompt, and carries noisy bubble artifacts, as shown in the second column in Fig. 7. We show that the multiresolution texture is capable of producing a visually appealing and highly detailed mesh appearance.

**Does cross-attention strengthen the style consistency?**



Figure 5. **Qualitative comparisons.** Latent-Paint [39] suffers from over-saturation and hallucinates scene components. MVDiffusion [63] delivers blurry textures and fails to reflect the input prompts. Text2Tex [8] struggles to keep all instances style-consistent. In contrast, our method produces high-quality textures and maintains overall style-consistency across instances in the scenes. Ceilings and back-facing walls are excluded for better visualizations. Images best viewed in color.

Replacing the cross-attention decoder module with a simple MLP also produces plausible textures. However, such replacement exposes global style inconsistency issue. Due to limited field of view and self-occlusion, the appearance of the same object can be synthesized differently. As shown in Fig. 7, big objects such as the carpet do not share a coherent pattern. It is difficult for the big objects to maintain style

consistency during optimization, if there is no global information shared across views. The proposed cross-attention decoder effectively tackles this issue by globally sharing the style features within each object. This enforces the instance style awareness, and therefore produces more style-consistent textures for all instances across the scene.





Figure 6. **Synthesized textures for 3D-FRONT scenes.** Our method generates high-quality style-coherent textures, and reflects the iconic traits in the prompts. Ceilings and back-facing walls are excluded for better visualizations. Images best viewed in color.



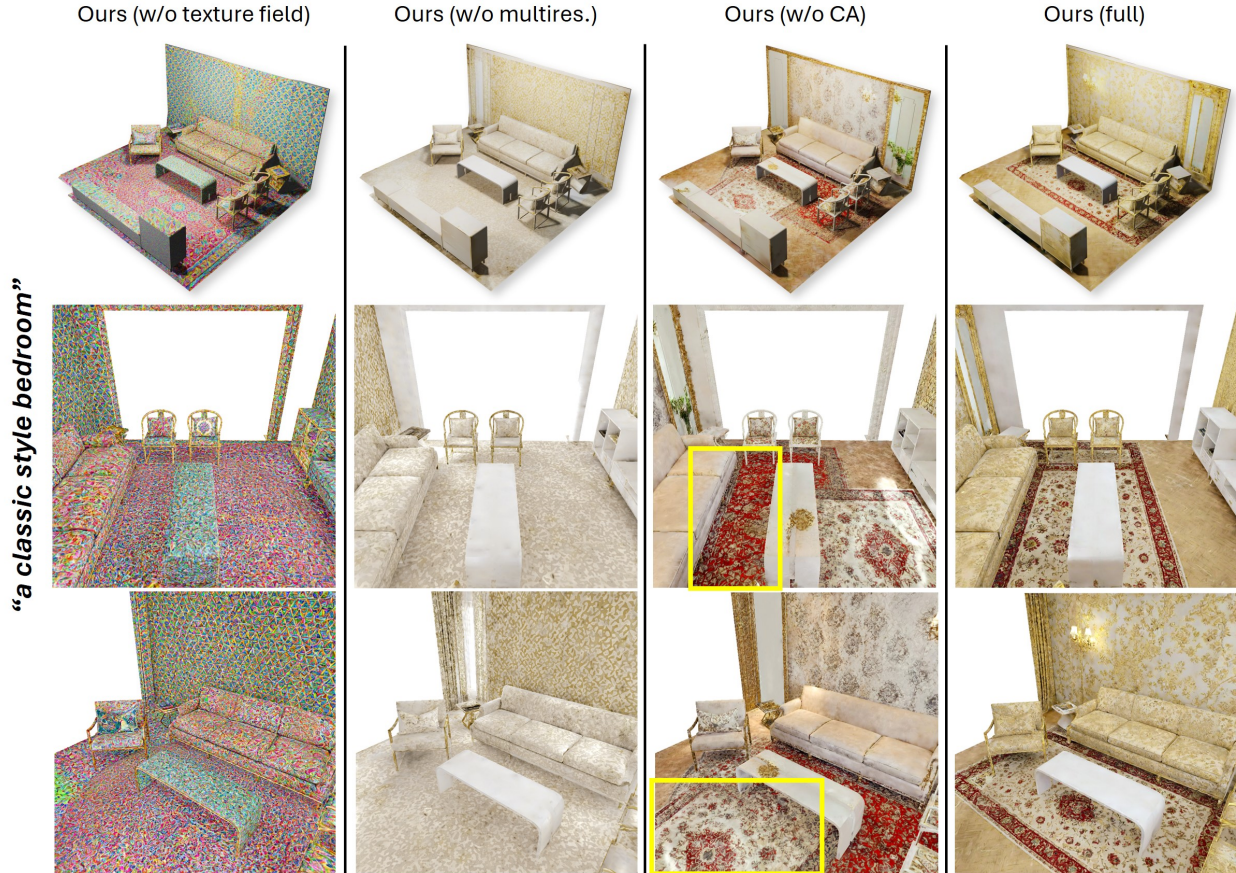


Figure 7. **Ablation studies on the key components.** Optimizing an RGB texture directly without the proposed texture field results in extreme noisy and unrealistic textures. A single-resolution texture fails to capture texture details and produces bubble artifacts. Removing the cross-attention decoder leads to style inconsistency, especially for big instances such as carpet, as shown in the yellow boxes. In contrast, our full method produces high-quality and style-consistent textures without aforementioned artifacts.

#### 4.5. Limitations

Although our method enables high-quality texture synthesis for indoor scenes, we still notice that our method tends to generate textures with shading effects. This phenomenon becomes more obvious when the scene structure indicates the existence of lighting such as lamp, window, or even mirror. We believe this issue can be properly addressed by carefully fine-tuning the diffusion priors on the indoor scene images without shading effects. We acknowledge this challenge and leave it to future research.

#### 5. Conclusion

We introduce SceneTex, a novel method for effectively generating high-quality and style-consistent textures for indoor scenes using depth-to-image diffusion priors. At its core, SceneTex proposes a multiresolution texture field to implicitly encode the mesh appearance. We optimize the target texture via a score-distillation-based objective function in respective RGB renderings. To further secure the style consistency across views, we introduce a cross-attention

decoder to predict the RGB values by cross-attending to the pre-sampled UV coordinates within each instance. We show that the proposed texture field with multiresolution texture is capable of generating visually appealing high-quality texture. Moreover, the proposed cross-attention decoder further strengthens the global style awareness for each instance, resulting in style-coherent appearance in the target scene. Extensive analysis show that SceneTex enables various and accurate texture synthesis for 3D-FRONT scenes, demonstrating significant improvements in visual quality and prompt fidelity over the prior texture generation methods. Overall, we hope our work can inspire more future research in the area of text-to-3D generation.

#### Acknowledgements

This work was supported by a gift by Snap Inc., the ERC Starting Grant Scan2CAD (804724), and the German Research Foundation (DFG) Research Unit “Learning and Simulation in Visual Computing”. We thank Angela Dai for the video voiceover.



## References

- [1] Rameen Abdal, Hsin-Ying Lee, Peihao Zhu, Menglei Chai, Aliaksandr Siarohin, Peter Wonka, and Sergey Tulyakov. 3davatarGAN: Bridging domains for personalized editable avatars. In *CVPR*, 2023. 2
- [2] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning representations and generative models for 3d point clouds. In *ICML*, 2018. 1, 2
- [3] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *CVPR*, 2021. 2
- [4] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *CVPR*, 2022. 2
- [5] Dave Zhenyu Chen, Angel X Chang, and Matthias Nießner. ScanRefer: 3D object localization in RGB-D scans using natural language. In *European Conference on Computer Vision*, pages 202–221. Springer, 2020. 2
- [6] Dave Zhenyu Chen, Ronghang Hu, Xinlei Chen, Matthias Nießner, and Angel X Chang. Unit3d: A unified transformer for 3d dense captioning and visual grounding. *arXiv preprint arXiv:2212.00836*, 2022.
- [7] Dave Zhenyu Chen, Qirui Wu, Matthias Nießner, and Angel X Chang. D 3 net: A unified speaker-listener architecture for 3d dense captioning and visual grounding. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXII*, pages 487–505. Springer, 2022. 2
- [8] Dave Zhenyu Chen, Yawar Siddiqui, Hsin-Ying Lee, Sergey Tulyakov, and Matthias Nießner. Text2tex: Text-driven texture synthesis via diffusion models. In *ICCV*, 2023. 1, 2, 5, 6
- [9] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. In *ICCV*, pages 22246–22256, 2023. 2
- [10] Wenzheng Chen, Huan Ling, Jun Gao, Edward Smith, Jaakko Lehtinen, Alec Jacobson, and Sanja Fidler. Learning to predict 3d objects with an interpolation-based differentiable renderer. 2019. 2
- [11] Yaosen Chen, Qi Yuan, Zhiqiang Li, Yuegen Liu, Wei Wang, Chaoping Xie, Xuming Wen, and Qien Yu. Upst-nerf: Universal photorealistic style transfer of neural radiance fields for 3d scene. *arXiv preprint arXiv:2208.07059*, 2022. 2
- [12] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *CVPR*, 2019. 1, 2
- [13] Zhenyu Chen, Ali Gholami, Matthias Nießner, and Angel X Chang. Scan2cap: Context-aware dense captioning in rgb-d scans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3193–3203, 2021. 2
- [14] Yen-Chi Cheng, Hsin-Ying Lee, Sergey Tulyakov, Alexander Schwing, and Liangyan Gui. Sdfusion: Multimodal 3d shape completion, reconstruction, and generation. In *CVPR*, 2023. 2
- [15] Zezhou Cheng, Menglei Chai, Jian Ren, Hsin-Ying Lee, Kyle Olszewski, Zeng Huang, Subhransu Maji, and Sergey Tulyakov. Cross-modal 3d shape generation and manipulation. In *ECCV*, 2022. 2
- [16] Pei-Ze Chiang, Meng-Shiun Tsai, Hung-Yu Tseng, Wei-Sheng Lai, and Wei-Chen Chiu. Stylizing 3d scene via implicit representation and hypernetwork. 2022. 2
- [17] Angela Dai, Yawar Siddiqui, Justus Thies, Julien Valentin, and Matthias Nießner. Spsg: Self-supervised photometric scene generation from rgb-d scans. In *CVPR*, 2021. 2
- [18] Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv preprint arXiv:2307.08691*, 2023. 5
- [19] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. 2021. 2
- [20] Huan Fu, Bowen Cai, Lin Gao, Ling-Xiao Zhang, Jiaming Wang, Cao Li, Qixun Zeng, Chengyue Sun, Rongfei Jia, Binqiang Zhao, et al. 3d-front: 3d furnished rooms with layouts and semantics. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10933–10942, 2021. 2, 5
- [21] Jun Gao, Tianchang Shen, Zian Wang, Wenzheng Chen, Kangxue Yin, Daiqing Li, Or Litany, Zan Gojcic, and Sanja Fidler. Get3d: A generative model of high quality 3d textured shapes learned from images. In *Advances In Neural Information Processing Systems*, 2022. 2
- [22] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *CVPR*, 2016. 2
- [23] Leon A Gatys, Alexander S Ecker, Matthias Bethge, Aaron Hertzmann, and Eli Shechtman. Controlling perceptual factors in neural style transfer. In *CVPR*, 2017. 2
- [24] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *NeurIPS*, 2014. 2
- [25] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d aware generator for high-resolution image synthesis. In *ICLR*, 2022. 2
- [26] Ayaan Haque, Matthew Tancik, Alexei A Efros, Alexander Holynski, and Angjoo Kanazawa. Instruct-nerf2nerf: Editing 3d scenes with instructions. *arXiv preprint arXiv:2303.12789*, 2023. 2
- [27] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. 2020. 2
- [28] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *arXiv preprint arXiv:2106.15282*, 2021. 2
- [29] Lukas Höllein, Justin Johnson, and Matthias Nießner. Stylemesh: Style transfer for indoor 3d scene reconstructions. In *CVPR*, 2022. 2
- [30] Lukas Höllein, Ang Cao, Andrew Owens, Justin Johnson, and Matthias Nießner. Text2room: Extracting textured 3d meshes from 2d text-to-image models. In *ICCV*, 2023. 2

- [31] Ronghang Hu and Amanpreet Singh. Unit: Multimodal multitask learning with a unified transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1439–1449, 2021. 2
- [32] Yi-Hua Huang, Yue He, Yu-Jie Yuan, Yu-Kun Lai, and Lin Gao. Stylizednerf: consistent 3d scene stylization as stylized nerf via 2d-3d mutual learning. In *CVPR*, 2022. 2
- [33] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 2
- [34] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 2
- [35] 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. *arXiv preprint arXiv:2211.10440*, 2022. 2
- [36] Chieh Hubert Lin, Hsin-Ying Lee, Willi Menapace, Menglei Chai, Aliaksandr Siarohin, Ming-Hsuan Yang, and Sergey Tulyakov. Infiniticity: Infinite-scale city synthesis. In *ICCV*, 2023. 1, 2
- [37] Shichen Liu, Tianye Li, Weikai Chen, and Hao Li. Soft rasterizer: A differentiable renderer for image-based 3d reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7708–7717, 2019. 3
- [38] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *CVPR*, 2021. 1, 2
- [39] Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. Latent-nerf for shape-guided generation of 3d shapes and textures. In *CVPR*, 2023. 1, 2, 3, 4, 5, 6
- [40] Paritosh Mittal, Yen-Chi Cheng, Maneesh Singh, and Shubham Tulsiani. AutoSDF: Shape priors for 3d completion, reconstruction and generation. In *CVPR*, 2022. 1, 2
- [41] Norman Müller, Yawar Siddiqui, Lorenzo Porzi, Samuel Rota Buló, Peter Kotschieder, and Matthias Nießner. Diffri: Rendering-guided 3d radiance field diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4328–4338, 2023. 2
- [42] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022. 5
- [43] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. 2021. 2
- [44] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *CVPR*, 2021. 2
- [45] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 5
- [46] Dario Pavllo, Graham Spinks, Thomas Hofmann, Marie-Francine Moens, and Aurelien Lucchi. Convolutional generation of textured 3d meshes. 2020. 2
- [47] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *ICLR*, 2023. 2, 3
- [48] 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 *ICLR*, 2021. 2, 5
- [49] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv:2007.08501*, 2020. 5
- [50] Elad Richardson, Gal Metzer, Yuval Alaluf, Raja Giryes, and Daniel Cohen-Or. Texture: Text-guided texturing of 3d shapes. *arXiv preprint arXiv:2302.01721*, 2023. 1, 2
- [51] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 4
- [52] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE TPAMI*, 2022. 2
- [53] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. 2020. 2
- [54] Aliaksandr Siarohin, Willi Menapace, Ivan Skorokhodov, Kyle Olszewski, Hsin-Ying Lee, Jian Ren, Menglei Chai, and Sergey Tulyakov. Unsupervised volumetric animation. In *CVPR*, 2023. 2
- [55] Aliaksandr Siarohin, Willi Menapace, Ivan Skorokhodov, Kyle Olszewski, Jian Ren, Hsin-Ying Lee, Menglei Chai, and Sergey Tulyakov. Unsupervised volumetric animation. In *CVPR*, 2023. 2
- [56] Yawar Siddiqui, Justus Thies, Fangchang Ma, Qi Shan, Matthias Nießner, and Angela Dai. Texturify: Generating textures on 3d shape surfaces. In *ECCV*, 2022. 2
- [57] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. Flava: A foundational language and vision alignment model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15638–15650, 2022. 2
- [58] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *Advances in neural information processing systems*, 33:7462–7473, 2020. 5
- [59] Ivan Skorokhodov, Aliaksandr Siarohin, Yinghao Xu, Jian Ren, Hsin-Ying Lee, Peter Wonka, and Sergey Tulyakov. 3d generation on imagenet. In *ICLR*, 2023. 2
- [60] Edward J Smith and David Meger. Improved adversarial systems for 3d object generation and reconstruction. In *Conference on Robot Learning*, pages 87–96. PMLR, 2017. 1, 2, 5
- [61] Gabriela Ben Melech Stan, Diana Wofk, Scottie Fox, Alex Redden, Will Saxton, Jean Yu, Estelle Aflalo, Shao-Yen Tseng, Fabio Nonato, Matthias Muller, et al. Ldm3d: Latent diffusion model for 3d. *arXiv preprint arXiv:2305.10853*, 2023. 2
- [62] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features



let networks learn high frequency functions in low dimensional domains. *Advances in Neural Information Processing Systems*, 33:7537–7547, 2020. 5

- [63] Shitao Tang, Fuyang Zhang, Jiacheng Chen, Peng Wang, and Yasutaka Furukawa. Mvdiffusion: Enabling holistic multi-view image generation with correspondence-aware diffusion. *arXiv preprint arXiv:2307.01097*, 2023. 5, 6
- [64] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. 2017. 2
- [65] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. In *CVPR*, 2023. 2
- [66] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. OFA: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning*, pages 23318–23340. PMLR, 2022. 2
- [67] 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. *arXiv preprint arXiv:2305.16213*, 2023. 2, 3, 4, 5
- [68] Jianwen Xie, Zilong Zheng, Ruiqi Gao, Wenguan Wang, Song-Chun Zhu, and Ying Nian Wu. Learning descriptor networks for 3d shape synthesis and analysis. In *CVPR*, 2018. 1, 2
- [69] Yinghao Xu, Menglei Chai, Zifan Shi, Sida Peng, Ivan Skokhodov, Aliaksandr Siarohin, Ceyuan Yang, Yujun Shen, Hsin-Ying Lee, Bolei Zhou, et al. Discoscene: Spatially disentangled generative radiance fields for controllable 3d-aware scene synthesis. In *CVPR*, 2023. 2
- [70] Rui Yu, Yue Dong, Pieter Peers, and Xin Tong. Learning texture generators for 3d shape collections from internet photo sets. In *BMVC*, 2021. 2
- [71] Kai Zhang, Nick Kolkin, Sai Bi, Fujun Luan, Zexiang Xu, Eli Shechtman, and Noah Snavely. Arf: Artistic radiance fields. In *ECCV*, 2022. 2
- [72] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023. 5
- [73] Song-Hai Zhang, Yuan-Chen Guo, and Qing-Wen Gu. Sketch2model: View-aware 3d modeling from single free-hand sketches. In *CVPR*, 2021. 1, 2
- [74] Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel diffusion. In *ICCV*, 2021. 2