
Emu3D: Text-to-Mesh Generation with High-Quality Geometry, Texture, and PBR Materials

Anonymous Author(s)

Affiliation

Address

email

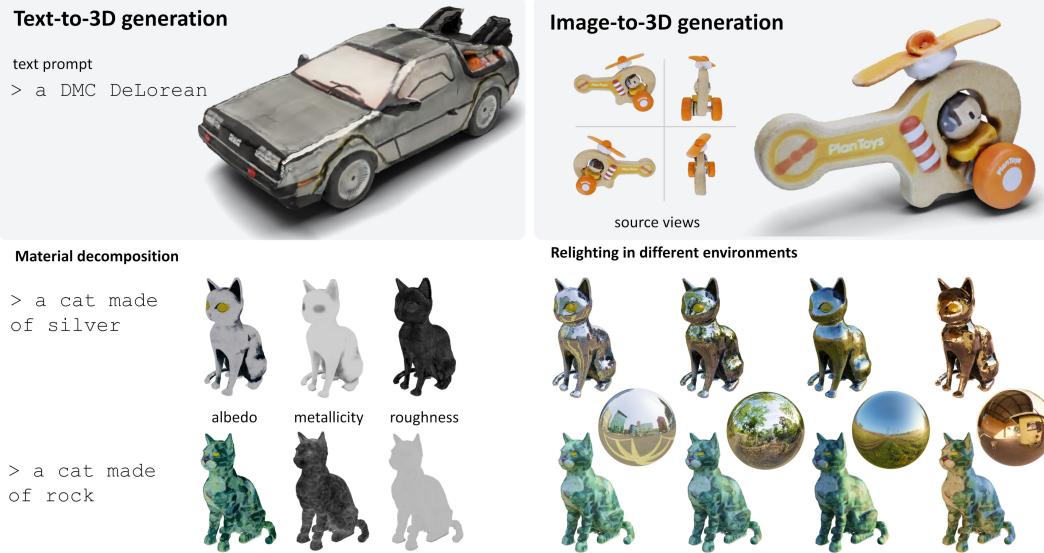


Figure 1: **Emu3D.** We present a novel text- or image-conditioned generator of 3D meshes with physically-based rendering materials (top). Emu3D produces meshes with detailed geometry, high-fidelity textures, and decomposes materials into albedo, metallicity, and roughness (bottom left), allowing us to realistically relight them in new environments (bottom right).

Abstract

1 3D generative AI has enormous potential for industry, but its quality remains
2 insufficient for professional use. We present Emu3D, a significant advancement
3 in text-to-3D which produces faithful, high-quality meshes with full material
4 control. Compared to prior works that bake shading in the 3D object, we output
5 physically-based rendering (PBR) materials, supporting realistic relighting. Emu3D
6 generates first several views of the object with factored shaded and albedo, and then
7 reconstructs colours, metallicity and roughness in 3D, training the reconstruction
8 network via a deferred shading loss. The latter models shape more reliably by
9 outputting directly a sign-distance function and using a corresponding loss, which,
10 implemented in fused kernels, are highly memory efficient. After mesh extraction, a
11 final texture refinement step operating in UV space significantly improves sharpness
12 and details. Emu3D achieves 17% improvement in Chamfer Distance and 40% in
13 LPIPS over the best concurrent work for few-view reconstruction, and a human
14 preference of 72% over the best industry text-to-3D tool, including those that
15 support PBR, with comparable inference time. Anonymized project page with
16 generated 3D assets: <https://anon-neurips-2024.github.io/>

17 **1 Introduction**

18 Generating 3D objects from text or images has enormous potential for 3D graphics, animation,
19 gaming and AR/VR. However, despite the outstanding progress of image and video generation [66,
20 41, 53, 39, 94, 100], the quality of 3D generators is still insufficient for professional use. Some
21 methods are excruciatingly slow and the generated 3D meshes and texture may contain defects.
22 Many approaches still “bake” appearance as albedo, ignoring the model’s response to environmental
23 lighting. This results in visually unattractive outputs, especially for reflective materials, which look
24 out of place when put in novel environments.

25 We introduce *Emu3D*, a significant step-up in text-conditioned 3D generation. Emu3D outperforms
26 prior works in faithfulness, quality of the generated 3D meshes and, especially, quality and control
27 of materials using *Physically-Based Rendering* (PBR) [84], while still generating assets in under
28 30 s. Emu3D uses a two-stage design from [39]. The first stage *stochastically* generates four images
29 of the object from canonical viewpoints, and the second stage *deterministically* reconstructs the
30 3D shape and appearance from these views (Fig. 1). This approach is faster and more robust than
31 SDS-based techniques [66] and produces more diverse and faithful results than single-stage 3D
32 generators [35, 59, 91, 76].

33 Emu3D explicitly models the interaction of light with an object’s surface using PBR, which is
34 essential for realistic 3D graphics. PBR in its simplest form accounts for albedo, metallicity, and
35 roughness to render scenes that accurately reflect environmental illumination. We focus on meshes as
36 the output representation due to their prevalence in applications and compatibility with PBR.

37 Our first contribution extends the two stages to generate PBR information. Since the image-to-3D
38 stage is deterministic, we task the stochastic text-to-image stage with resolving material-assignment
39 ambiguities. Instead of generating textures with baked light, we aim to fine-tune a text-to-image
40 model to output partial PBR channels. However, since the distribution of PBR images differs from
41 natural images, direct fine-tuning is challenging. Our solution is to generate both shaded and albedo
42 (unshaded) channels in the 4-view stage. This allows the image-to-3D component to accurately
43 predict PBR materials by analyzing the differences between the albedo and shaded textures.

44 Our second innovation is in the mesh generation step. Existing works typically output an implicit
45 shape representation but use an opacity field, which can lead to ill-defined level sets, which result
46 in meshing artifacts. To address this, Emu3D’s image-to-3D stage, *EmuLRM*, directly predicts a
47 signed-distance field (SDF) instead of an opacity field. This yields higher-quality meshes, as the
48 zero-level-set of an SDF accurately traces the object’s surface and can be directly supervised with
49 ground-truth depth maps, which is not immediately possible for opacities. We base our SDF predictor
50 on VolSDF [106] while improving training efficiency by implementing rendering within the memory-
51 efficient Lightplane kernels [6]. This allows for larger batches and photometric loss supervision on
52 high-resolution renders, leading to better texture quality.

53 In the mesh generation step, we also add a novel texture refiner that learns to backproject the input
54 view’s texture onto the extracted albedo and materials, resolving potential conflicts from the different
55 views. The refiner significantly sharpens the texture channels, whose high-frequency details can
56 otherwise get washed out by the image-to-mesh reconstructor and mesh conversion.

57 We demonstrate the effectiveness of Emu3D on both image-to-3D and text-to-3D tasks. For image-to-
58 3D, we attain state-of-the-art performances among existing few-view mesh-reconstruction methods
59 when measuring the accuracy of both shaded and PBR texture maps. For text-to-3D, we conduct
60 extensive user studies comparing to the best academic methods and industry tools with comparable
61 inference time, outperforming them all in terms of visual quality and text alignment.

62 **2 Related Work**

63 **Text-to-3D.** Inspired by text-to-image models, early approaches [62, 32, 25, 108, 103] train 3D
64 diffusion models on datasets of captioned 3D assets. Yet, the limited size and diversity of 3D data
65 prevents generalization to open-vocabulary prompts. Recent works thus leverage text-to-image
66 models trained on billion-scale text-annotated image data. Works like [73, 54] propose to directly
67 finetune 2D diffusion models to output 3D representations, but the quality is typically limited because
68 of the high domain gap. The rest of text-to-3D works is two-fold.

69 The first are distillation models building on the seminal DreamFusion [66], where SDS loss optimizes
70 a NeRF to produce renders that match the belief of a pre-trained text-to-image model. Extensions
71 study: (i) other 3D representations like hash grids [41, 67], meshes [41] or 3D gaussians [78, 109, 12];
72 (ii) improved SDS [88, 92, 117, 29]; (iii) monocular conditioning facilitating realism [67, 79, 111, 75];
73 (iv) predicting additional normals or depth for better geometry [68, 75]. Yet, distillation methods are
74 prone to 3D abnormalities like the multi-face (Janus) problem or 3D content drift [72]. A common
75 solution is to incorporate view-consistency priors into the diffusion model, by either conditioning on
76 cameras [44, 71, 30, 11, 67] or by generating multiple object views [72, 95, 90, 37, 116]. Additionally,
77 slow SDS optimization restricts the generation time to several minutes per asset, which is partly
78 addressed in [50, 98] with amortized SDS.

79 The second category comprises faster two-stage approaches that reconstruct multiple object views
80 generated with a text-to-image model. [43, 48, 46, 105, 104, 8, 80, 27, 22] fine-tuned text-to-
81 image/video models [52, 13] to output accurate multi-view generations followed by per-scene
82 optimization like NeRF [56] or 3DGs [36]. However, this requires many perfectly consistent views
83 which are difficult to generate with current text-to-image models. Instant3D [39] improves runtime
84 by generating a 4-view image grid followed by a feed-forward LRM [28] 4-view reconstructor.
85 One-2-3-45++ [42] replaces the LRM with a 3D diffusion model. In this work, we adopt Instant3D
86 paradigm and, in addition to LRM changes enabling improved 3D quality and support for material
87 prediction, we extend the 4-view generation part to predict additional material information. We found
88 that this modification plays an important role to predict accurate 3D shape and materials from images.

89 **3D reconstruction from images.** 3D scene reconstruction, in its traditional *multi-view stereo* (MVS)
90 form, assumes access to a dense set of scene views. Recent methods, such as NeRF [56], optimize a
91 3D representation by minimizing multi-view rendering errors. The underlying 3D representations are
92 two-fold: (i) explicit representations like meshes [21, 112, 23, 61, 57] or 3D points / gaussians [36, 24],
93 and (ii) implicit neural representations like occupancy fields [63], radiance fields [56, 60] or signed
94 distance functions (SDF) [107]. Compared to occupancy fields, SDF [64, 106, 89, 17, 20] simplifies
95 surface constraints integration, improving scene geometry. For this reason, we adopt an SDF
96 formulation and demonstrate that it outperforms occupancy.

97 *Sparse-view reconstruction* instead assumes few input views (usually 1 to 8). This under-constrained
98 problem is often tackled by training feed-forward reconstructors on large datasets [14, 34, 55, 45,
99 97, 58, 91]. In particular, the state-of-the-art LRM [28] trains a large Transformer [86] to predict a
100 triplane-supported NeRF [7, 9]. LRM extensions study other 3D representations like meshes [100, 94]
101 or 3D gaussians [118, 102, 77, 113], improved backbones [93, 94] and training protocols [83, 33].
102 Our proposed approach also builds upon LRM but introduces three key modifications: (i) an SDF
103 formulation for improved geometry, (ii) a PBR appearance modeling for material estimation, (iii) a
104 texture refiner for better texture details. Note that a parallel effort [53, 96] aims mitigating the need
105 for dense multi-views in per-scene optimization frameworks by leveraging 2D diffusion priors.

106 **3D modeling with PBR appearance.** Most 3D modelling methods output 3D shapes with emulated
107 light interaction [56, 36] or with baked surface lighting [28]. Since baked lighting ignores model's
108 response to environmental light, it is unsuitable for graphics pipelines with controlled lighting.
109 Physically-based rendering (PBR) defines material properties so that a suitable shader outputs
110 physically plausible renders. As such, works in MVS develop PBR-extensions for existing 3D
111 representations such as NeRF [4, 3, 99], SDF modeling [114], differentiable meshes [61, 26] or
112 3D gaussians [31, 40]. In generative modelling, [10, 68, 47, 101] augment the text-to-3D SDS
113 optimization [66] with a PBR model. Differently from them, we integrate PBR to our feed-forward
114 text-to-3D model, unlocking for the first time the fast text-based generation of PBR assets.

115 3 Method

116 Emu3D is a two-stage pipeline (Fig. 2). First, text-to-image (Sec. 3.1), takes text as input and
117 generates a 4-view grid of images with material information. Second, image-to-3D, comprises a
118 novel PBR-based sparse-view reconstruction model (Sec. 3.2) and the texture refiner (Sec. 3.3). As
119 such, Emu3D is applicable to two tasks: text-to-3D (stage 1+2) and image-to-3D (stage 2 only).

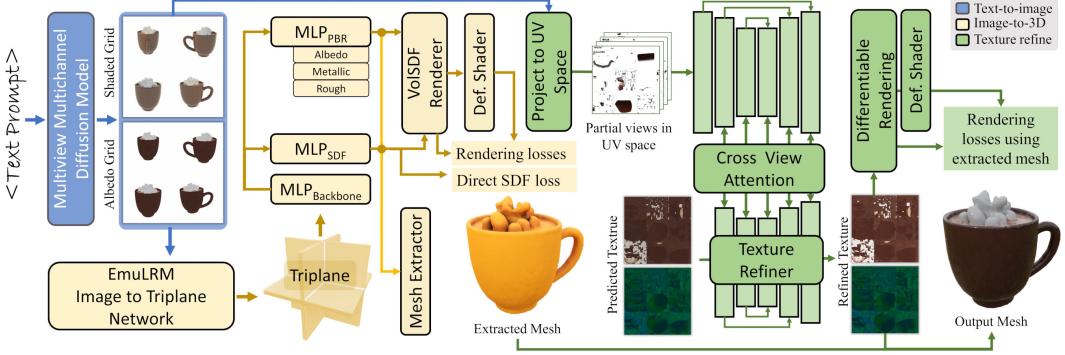


Figure 2: **Overview.** Given a text prompt, Emu3D generates a 3D mesh with PBR materials in two stages. The first text-to-image stage (blue) predicts a 6-channel image depicting 4 views of the object with shaded and albedo colors. The second image-to-3D stage includes two steps. First, a 3D reconstructor (dubbed EmuLRM) outputs a triplane-supported SDF field converted into a mesh with textured PBR materials (orange). Then, PBR materials are enhanced with our texture refiner which recovers missing details from the input views (green).

120 3.1 Text-to-image: Generating shaded and albedo images from text

121 The goal of the text-to-image module is to generate several views of the generated 3D object. To this
 122 end, we employ an internal text-to-image diffusion model pre-trained on billions of text-annotated
 123 images, with an architecture similar to Emu [16]. Similar to [72, 39], we finetune the model to
 124 predict a grid of four images $I_i, i = 1, \dots, 4$, each depicting the object from canonical viewpoints π_i .
 125 Note that I_i are RGB images of the shaded object. In principle, we could defer the PBR parameter
 126 extraction to the image-to-3D stage, but this led to suboptimal results. This is due to the determinism
 127 of the image-to-3D stage which fails to model ambiguities when assigning materials to surfaces.

128 A natural solution, then, is to predict the PBR parameters directly in the text-to-image stage. These
 129 consists of the albedo ρ_0 (by which we mean the base color, which technically is different for non-zero
 130 metallicity), the metallicity γ , and the roughness α . However, we found this to be ineffective too
 131 because the metallicity and roughness maps deviate from the distribution of natural images making
 132 them a hard target for finetuning. Our novel solution is to train the model to generate instead a **4-view**
 133 **grid with 6 channels**, 3 for the shaded appearance I and 3 more for the albedo ρ_0 . This reduces the
 134 finetuning gap, and removes enough ambiguity for accurate PBR prediction in the image-to-3D stage.

135 3.2 Image-to-3D: A PBR-based large reconstruction model

136 We now describe the image-to-3D stage, which solves the reconstruction tasks given either a small
 137 number of views I_i (few-view reconstruction), or the 4-view 6-channel grid of Sec. 3.1.

138 At the core of our method is a new PBR-aware reconstruction model, EmuLRM, that reconstructs
 139 the object given $N \geq 1$ *posed* images $(I_i, \pi_i)_{i=1}^N$, where $I_i \in \mathbb{R}^{H \times W \times D}$ and $\pi_i \in \Pi$ is the camera
 140 viewpoint. As noted in Sec. 3.1, we consider $N = 4$ canonical viewpoints π_1, \dots, π_4 (fixed to
 141 20° elevation and $0^\circ, 90^\circ, 180^\circ, 270^\circ$ azimuths) and $D = 6$ input channels. The output is a 3D
 142 field representing the shape and PBR materials of the object as an SDF $s : \mathbb{R}^3 \rightarrow \mathbb{R}$, where $s(\mathbf{x})$
 143 is the signed distance from the 3D point \mathbf{x} to the nearest object surface point, and a PBR function
 144 $k : \mathbb{R}^3 \rightarrow \mathbb{R}^5$, where $k(\mathbf{x}) = (\rho_0, \gamma, \alpha)$ are the albedo, metallicity and roughness.

145 The key to learning the model is the *differentiable rendering* operator \mathcal{R} . This takes as input a field
 146 $\ell : \mathbb{R}^3 \rightarrow \mathbb{R}^D$, the SDF s , the viewpoint π , and a pixel $u \in U = [0, W] \times [0, H]$, and outputs the
 147 projection of the field on the pixel according to the rendering equation [56], which has the same
 148 number of channels D as the rendered field ℓ :

$$\mathcal{R}(u | \ell, s, \pi) = \int_0^\infty \ell(\mathbf{x}_t) \sigma(\mathbf{x}_t | s) e^{-\int_0^t \sigma(\mathbf{x}_\tau | s) d\tau} dt. \quad (1)$$

149 Here $\mathbf{x}_t = \mathbf{x}_0 - t\omega_0$, $t \in [0, \infty)$ is the ray that goes from the camera center \mathbf{x}_0 through the pixel u
 150 along direction $-\omega_0 \in \mathbb{S}^2$. The function $\sigma(\mathbf{x} | s)$ is the opacity of the 3D point \mathbf{x} and is obtained

151 from the SDF value $s(\mathbf{x})$ using the VolSDF [106] formula

$$\sigma(\mathbf{x} | s) = \frac{a}{2} \left(1 + \text{sign } s(\mathbf{x}) (1 - e^{-|s(\mathbf{x})|/b}) \right), \quad (2)$$

152 where a, b are the hyperparameters. We use Eq. (1) to render several different types of fields ℓ , the
153 most important of which is the radiance field, introduced next along with the material model.

154 **Reflectance model.** The appearance of the object $\hat{I}(u) = \mathcal{R}(u | L, s, \pi)$ in a shaded RGB image
155 \hat{I} is obtained by rendering its *radiance field* $\ell(\mathbf{x}) = L(\mathbf{x}, \omega_o | k, \mathbf{n})$, where \mathbf{n} is the field of unit
156 normals. The radiance is the light reflected by the object in the direction ω_o of the observer (see
157 App. A.8 for details). Key to PBR is that the radiance is the result of light reflected by the object:

$$L(\mathbf{x}, \omega_o | k, \mathbf{n}) = \int_{H(\mathbf{n})} f(\omega_i, \omega_o | k(\mathbf{x}), \mathbf{n}(\mathbf{x})) L(\mathbf{x}, -\omega_i)(\mathbf{n}(\mathbf{x}) \cdot \omega_i) d\Omega_i, \quad (3)$$

158 where $\omega_o, \omega_i \in H(\mathbf{n}) = \{\omega \in \mathbb{S}^2 : \mathbf{n} \cdot \omega \geq 0\}$ are two unit vectors pointing outside the object and
159 $L(\mathbf{x}, -\omega_i)$ is the radiance incoming from the environment at \mathbf{x} from direction ω_i in the solid angle
160 $d\Omega_i$. The *Bidirectional Reflectance Distribution Function* (BRDF) f tells how light received from
161 direction $-\omega_i$ (incoming) is scattered into different directions ω_o (outgoing) by the object [19].

162 In PBR, we consider a physically-inspired model for the BRDF, striking a balance between realism
163 and complexity [2, 84, 74, 15, 87]; specifically, we use the GGX model [87, 5], which depends on
164 parameters ρ_0 , γ , and α only (see App. A.12.1 for the parametric form of f). Hence, the EmuLRM
165 predicts the triplet $k(\mathbf{x}) = (\rho_0, \gamma, \alpha)$ at each 3D point \mathbf{x} .

166 **Deferred shading.** In practice, instead of computing $\hat{I}(u) = \mathcal{R}(u | L, s, \pi)$ using Eqs. (1) and (3),
167 we use the process of *deferred shading* [19]:

$$\hat{I}(u) = \mathcal{R}_{\text{def}}(u | k, s, \pi) = \int_{H(\mathbf{n})} f(\omega_i, \omega_o | \bar{k}, \bar{\mathbf{n}}) L_{\text{env}}(-\omega_i)(\bar{\mathbf{n}} \cdot \omega_i) d\Omega_i, \quad (4)$$

168 where L_{env} is the environment radiance (assumed to be the same for all \mathbf{x}), $\bar{k} = \mathcal{R}(u | k, s, \pi)$
169 and $\bar{\mathbf{n}} = \mathcal{R}(u | \mathbf{n}, s, \pi)$ are rendered versions of the material and normal fields. The advantage
170 of Eq. (4) is that the BRDF f is evaluated only once per pixel, which is much faster and less
171 memory intensive than doing so for each 3D point during the evaluation of Eq. (1), particularly for
172 training/backpropagation. During training, furthermore, the environment light is assumed to be a
173 single light source at infinity, so the integral (4) reduces to evaluating a single term.

174 **Training formulation and losses.** EmuLRM is thus a neural network that takes as input a set of
175 images $(I_i, \pi_i)_{i=1}^N$ and outputs estimates \hat{s} and \hat{k} for the SDF and PBR fields. We train it from a
176 dataset of mesh surfaces $M \subset \mathbb{R}^3$ with ground truth PBR materials $k : M \rightarrow \mathbb{R}^5$.

177 Reconstruction models are typically trained via supervision on renders [28, 94]. However, physically
178 accurate rendering via Eq. (1) is very expensive. We overcome this hurdle in two ways. First, we
179 render the raw ground-truth PBR fields k and use them to supervise their predicted counterparts with
180 the MSE loss, skipping Eq. (1). For the rendered albedo ρ_0 – which is similar enough to natural
181 images – we also use the LPIPS [115] loss:

$$\mathcal{L}_{\text{pbr}} = \text{LPIPS}\left(\mathcal{R}(\cdot | \hat{\rho}_0, \hat{s}, \pi), \mathcal{R}(\cdot | \rho_0, M, \pi)\right) + \left\| \mathcal{R}(\cdot | \hat{k}, \hat{s}, \pi) - \mathcal{R}(\cdot | k, M, \pi) \right\|^2. \quad (5)$$

182 We further supervise the PBR field by adding a computationally-efficient deferred shading loss:

$$\mathcal{L}_{\text{def}} = \|\sqrt{w} \odot (\mathcal{R}_{\text{def}}(\cdot | \hat{k}, \hat{s}, \pi) - I)\|^2, \quad (6)$$

183 where I is the ground-truth rendered image of the object. The weight $w(u) = \hat{\mathbf{n}}(u) \cdot \mathbf{n}(u)$ is the dot
184 product of the predicted and ground-truth normals at pixel u . It discounts the loss where the predicted
185 geometry is not yet learnt. Fig. 13 (b) visualizes deferred shading and the rendering loss.

186 Finally, we also supervise the SDF field with a direct loss \mathcal{L}_{sdf} (implemented as in [1]), a depth-MSE
187 loss $\mathcal{L}_{\text{depth}}$ between the depth renders and the ground truth, and with a binary cross-entropy $\mathcal{L}_{\text{mask}}$
188 between the alpha-mask renders and the ground-truth masks. Refer to App. A.6.2 for more details.

189 **LightPlane implementation.** We base EmuLRM on LightplaneLRM [6], a variant of LRM [28]
190 exploiting memory and compute-efficient Lightplane splatting and rendering kernels, offering better
191 quality reconstructions. However, since LightplaneLRM uses density fields, which are suboptimal for
192 mesh conversion [89, 64, 1], we extend the Lightplane rendering GPU kernel with a VolSDF [106]
193 renderer using Eq. (2). Additionally, we also fuse into the kernel the direct SDF loss \mathcal{L}_{sdf} since a
194 naive autograd implementation is too memory-heavy.

195 **3.3 Mesh extraction and texture refiner**

196 The EmuLRM module of Sec. 3.2 outputs a sign distance function s , implicitly defining the object
 197 surface $A = \{x \in \mathbb{R}^3 \mid s(x) = 0\}$ as a level set of s . We use the Marching Cubes algorithm [49] to
 198 trace the level set and output a mesh $M \approx A$. Then, xatlas [110] extracts a UV map $\phi : [0, V]^2 \rightarrow M$,
 199 mapping each 2D UV-space point $v = \phi(x)$ to a point $x \in M$ on the mesh.

200 Next, the goal is to extract a high-quality 5-channel PBR texture image $\bar{K} \in \mathbb{R}^{V \times V \times 5}$ capturing the
 201 albedo, metallicity, and roughness of each mesh point. The texture image K can be defined directly
 202 by sampling the predicted PBR field \hat{k} as $K(v) \leftarrow \hat{k}(\phi(v))$, but this often yields blurry results due to
 203 the limited resolution of EmuLRM. Instead, we design a texture refiner module which takes as input
 204 the coarse PBR-sampled texture image as well as the N views representing the object and outputs a
 205 much sharper texture \bar{K} . In essence, this module leverages the information from the different views
 206 to refine the coarse texture image. The right part of Fig. 2 illustrates such a module.

207 More specifically, it relies on a network Φ which is fed with $N + 1$ texture images $\{K_i\}_{i=0}^N$. First,
 208 each pixel $v \in [0, V]^2$ of $K_0 \in \mathbb{R}^{V \times V \times 11}$ is annotated with the concatenation of the normal, the
 209 3D location, and the output of EmuLRM’s PBR field $k(\phi(v))$ evaluated at v ’s 3D point $\phi(v)$. The
 210 remaining K_1, \dots, K_N correspond to partial texture images with 6 channels (for the base and shaded
 211 colors) which are obtained by backprojecting the object views to the mesh surface. The network Φ
 212 utilises two U-Nets to fuse $\{K_i\}_{i=0}^N$ into the enhanced texture \bar{K} . Φ ’s goal is to select, for each UV
 213 point v , which of the N input views provides the best information. Specifically, each partial texture
 214 image K_i is processed in parallel by a first U-Net, and the resulting information is communicated via
 215 cross attention to a second U-Net whose goal is to refine K_0 into the enhanced texture \bar{K} . Please
 216 refer to App. A.7 for further details.

217 Such a network is trained on the same dataset and supervised with the PBR and albedo rendering losses
 218 as EmuLRM. The only difference is meshes (whose geometry is fixed) are rendered differentiably
 219 using PyTorch3D’s [69] mesh rasterizer instead of the Lightplane SDF renderer.

220 **4 Experiments**

221 Our **training data** consists of 140,000 meshes of diverse semantic categories created by 3D artists.
 222 For each asset, we render 36 views at random elevations within the range of $[-30^\circ, 50^\circ]$ at uniform
 223 intervals of 30° around the object, lit with a randomly selected environment map. We render
 224 the shaded images, albedo, metallicity, roughness, depth maps, and foreground mask from each
 225 viewpoint. The text-to-image stage is based on an internal text-to-image model architecturally similar
 226 to Emu [16], fine-tuned on a subset of 10,000 high-quality 3D samples, captioned by a Cap3D-like
 227 pipeline [51] that uses Llama3 [85]. The other stage utilizes the entire 3D dataset instead.

228 For **evaluation**, following [102, 100, 39], we assess visual quality using PSNR and LPIPS [115]
 229 between the rendered and ground-truth images. PSNR is computed in the foreground region to avoid
 230 metric inflation due to the empty background. Geometric quality is measured by the L1 error between
 231 the rendered and ground-truth depth maps (of the foreground pixels), as well as the IoU of the object
 232 silhouette. We further report Chamfer Distance (CD) and Normal Correctness (NC) for 20,000

Table 1: **Four-view reconstruction with PBR** evaluating the accuracy of the PBR renders for EmuLRM and ablations. Methods in top / bottom accept 4 views with shaded / shaded&albedo color channels.

Method	LPIPS↓		PSNR↑		
	albedo	albedo	metal	rough	
C = LightplaneLRM w/ SDF	0.117	17.14	12.39	15.25	
E = C + Material prediction	0.097	20.66	15.99	20.25	
F = E + Deferred shading loss	0.093	21.12	18.64	20.66	
G = F + Texture refinement	0.087	21.97	22.19	20.85	
H = F + Albedo & shaded input	0.084	23.02	20.43	21.18	
I = H + Texture refinement	0.069	24.39	27.28	20.63	

Table 2: **Win-rate of Emu3D in text-to-3D user study** evaluating visual quality and the alignment between the prompt and the generated meshes. Emu3D beats all baselines at 30 sec budget (on an A100 GPU).

Method	Visual quality	Text fidelity	PBR
GRM [102]	96.7 %	93.3 %	✗
InstantMesh [100]	99.3 %	97.3 %	✗
LightplaneLRM [6]	66.6 %	N/A	✗
Meshy v3 [82]	94.6 %	91.3 %	✓
Luma Genie 1.0 [81]	72.3 %	72.8 %	✓

233 points uniformly sampled on both the predicted and ground-truth shapes. Material decomposition
 234 is evaluated with LPIPS and PSNR on the albedo image, and PSNR alone for the metallicity and
 235 roughness channels. All metrics are calculated on meshified outputs rather than on neural renders.

236 4.1 Sparse-view reconstruction

237 We tackle the sparse-view reconstruction
 238 task of predicting a 3D mesh
 239 from 4 posed images of an object
 240 on a subset of 332 meshes from
 241 Google Scanned Objects (GSO) [18].
 242 We compare against state-of-the-art
 243 Instant3D-LRM [39], GRM [102], In-
 244 stantMesh [100], and MeshLRM [94].
 245 We also include LightplaneLRM [6],
 246 an improved version of Instant3D-
 247 LRM, which serves as our base model.
 248 MeshLRM [94] has not been open-
 249 sourced so we compare only qualita-
 250 tively to meshes from their webpage.
 251 All methods are evaluated using the
 252 same input views at 512^2 resolution. Since none of the latter predict PBR materials and since
 253 GSO lacks ground-truth PBR materials, for fairness, we use a variant of our model that predicts
 254 shaded object textures.

255 As shown in Figs. 4 and 8 and Tab. 3, our method outperforms all baselines across all metrics. GRM
 256 captures texture detail well but struggles with fine geometric structures when meshified. InstantMesh
 257 and LightplaneLRM improve geometry but fall short on finer details and texture quality. Our approach
 258 excels in reconstructing shapes with detailed geometry and high-fidelity textures.

259 Ablations in Tab. 3 and Fig. 4 show that incorporating
 260 our scalable SDF-based rendering and direct SDF
 261 loss into the base LightplaneLRM model enhances
 262 geometric quality. Adding texture refinement further
 263 brings fine texture details; furthermore, as shown in
 264 Fig. 14, cross-view texture attention in UV space is
 265 crucial for performance.

266 Next, we consider the task of **sparse-view recon-**
267 struction with PBR materials, where the goal is to
 268 reconstruct the 3D geometry and texture properties
 269 (albedo, metallicity, and roughness) from 4 posed
 270 shaded 2D views of an object. This is done on an
 271 internal dataset of 256 artist-created 3D meshes, cu-
 272 rated for high-quality materials. Since there are no
 273 existing few-view feed-forward PBR reconstructors,
 274 we conduct an ablation study in Tab. 1 and Figs. 3
 275 and 12. While adding material prediction with addi-
 276 tional MLP heads provides some improvements, we
 277 observe that incorporating the deferred shading loss
 278 and texture refinement is essential for high-quality
 279 PBR decomposition. Example PBR predictions are
 280 shown in Fig. 7.

281 4.2 Text-to-3D generation

282 Finally, we evaluate text-to-3D with PBR materials. We compare against state-of-the-art feed-forward
 283 methods that generate assets at comparable speed (≈ 10 to 30 sec per asset). This includes text-to-3D
 284 variants of GRM [102], InstantMesh [100], and LightplaneLRM [6]. GRM uses Instant3D’s 4-view
 285 grid generator, InstantMesh receives the first view from our 2D diffusion model and subsequently

Table 3: **Four-view reconstruction on GSO** comparing the appearance and geometry of EmuLRM (outputting baked-light texture) to baselines (top) and ablations (bottom). CD values multiplied by 10^{-2} .

Method	LPIPS \downarrow	PSNR \uparrow	Depth \downarrow	IoU \uparrow	CD \downarrow	NC \uparrow
Instant3D-LRM [28]	0.124	18.54	0.325	0.930	1.630	0.844
GRM [102]	0.100	19.87	0.364	0.949	1.490	0.873
InstantMesh [100]	0.113	20.63	0.334	0.937	1.364	0.848
EmuLRM (ours)	0.057	22.49	0.173	0.968	1.137	0.885
A = LightplaneLRM [6]	0.095	18.60	0.456	0.953	1.313	0.872
B = A + VolsDF rendering	0.094	20.91	0.201	0.957	1.212	0.875
C = B + Direct SDF loss	0.083	21.75	0.173	0.968	1.137	0.885
D = C + Texture refinement	0.057	22.49	0.173	0.968	1.137	0.885

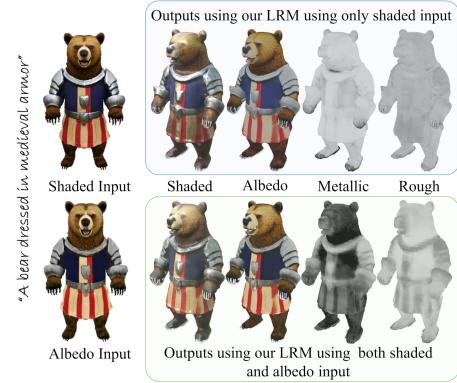


Figure 3: **Qualitative ablation on albedo generation.** In text-to-3D, generating 4 views representing albedo colors alongside shaded RGB colors improves material estimation for our 3D reconstructor. With both inputs, the model accurately predicts the armor as metallic and smooth, while the bear’s fur is rough.

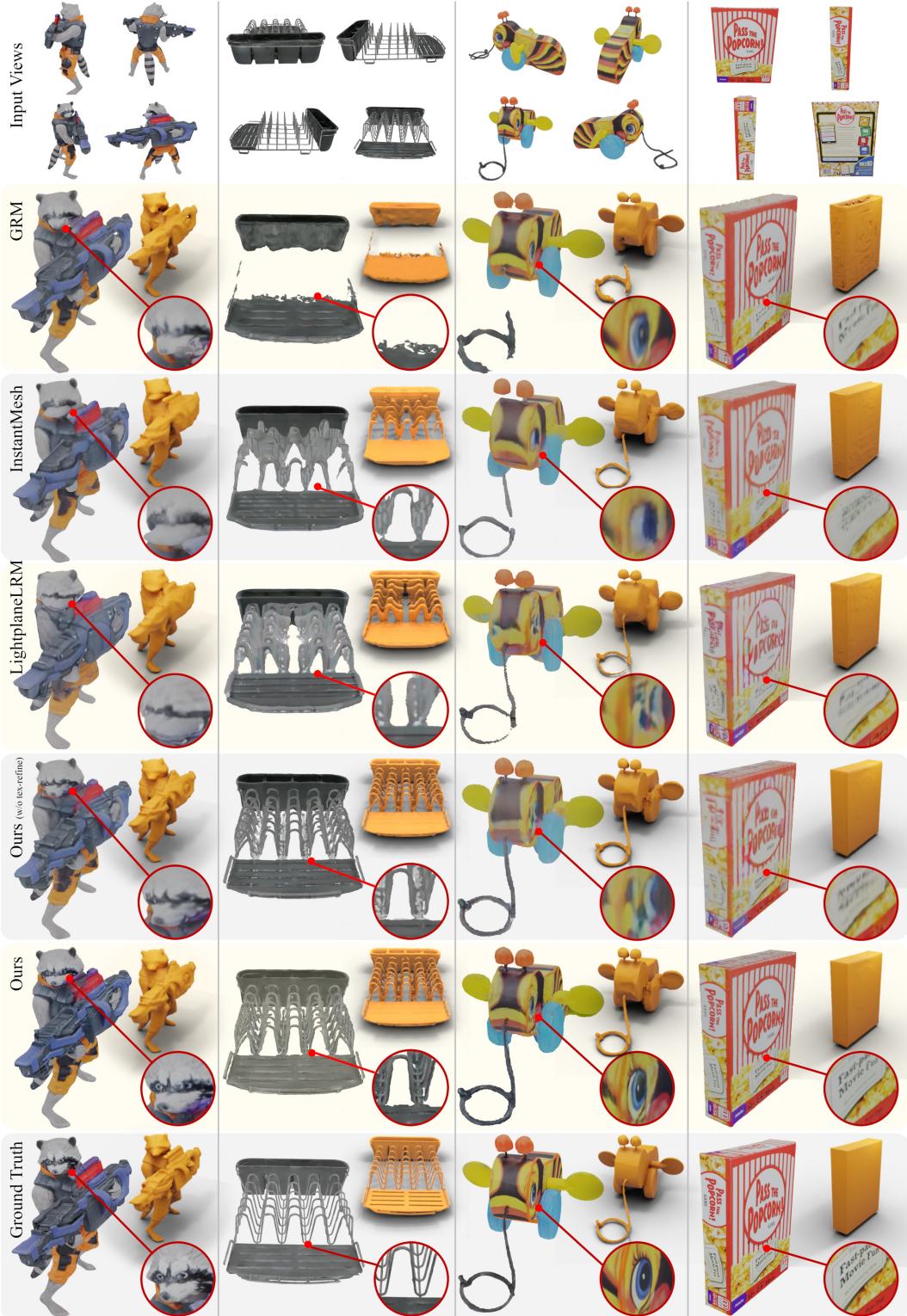


Figure 4: Qualitative comparison for sparse-view reconstruction. Our Emu3D gives better geometry (shown in orange) and higher fidelity texture (inset) compared to state of the art. SDF representation along with the direct SDF loss gives a better geometry compared to the base LightplaneLRM model which uses occupancy (row 4 and 5). Furthermore, our texture refiner greatly enhances texture fidelity (row 5 and 6).

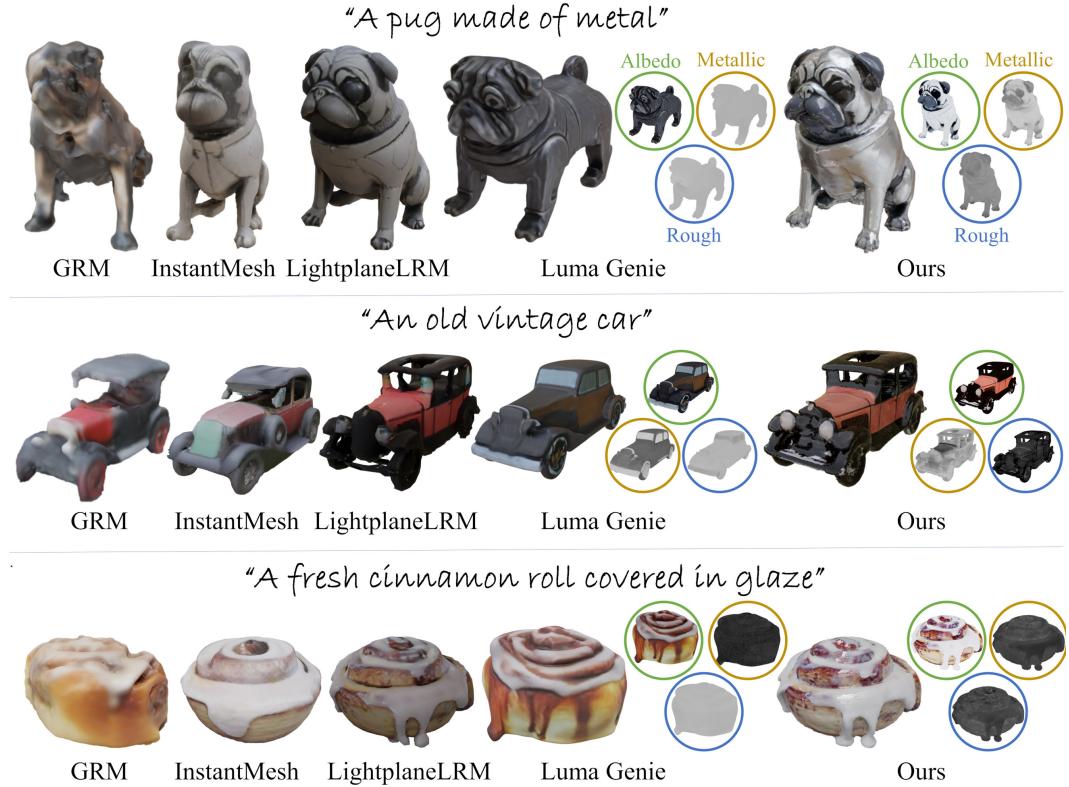


Figure 5: **Qualitative comparison for text-to-3D.** We compare 3D meshes generated by Emu3D and state-of-the-art baselines. We include material decomposition for methods producing PBR materials (Luma Genie and our Emu3D). Our approach produces higher quality materials with better-defined metallicity and roughness, and a more accurate decoupling of lighting effects in the albedo.

generates 6 views, while LightplaneLRM accepts 4 views from our grid generator. Since these methods bake lighting instead of generating PBR materials, for evaluation we apply flat texture shading to our outputs. Additionally, we compare with Meshy-3 [82] and LumaAI Genie [81], proprietary text-to-3D methods with PBR workflow capable of creating assets within 30 and 15 sec respectively. Fig. 5 shows that Emu3D meshes are visually more appealing and have meaningful materials. Figs. 9 and 11 provide more examples and showcase fine-grained material control.

For quantitative evaluation, we conducted an extensive user study in Tab. 2 using the 404 deduplicated text prompts from DreamFusion [66]. Users were shown 360° videos of the generated and baseline meshes and were asked to rate them based on 3D shape quality and alignment with the text prompt. A total of 11,080 responses were collected, with significant preference for Emu3D’s meshes.

Finally, we ablate the effect of generating dual-channel albedo+shaded grids compared to albedo-only input in Fig. 3 revealing significant PBR decomposition superiority of the former. Additionally, Fig. 12 illustrates the effect of our deferred shading loss.

5 Conclusions

We have introduced Emu3D, a significant advancement in sparse-view reconstruction and text-to-3D. Emu3D can generate 3D meshes with high-quality textures and PBR materials faithful to the input text. This uses several key innovations: generating multi-view grids with both shaded and albedo channels, introducing a new reconstruction network that predicts PBR materials from this information, using deferred shading to train this network, improving geometry via a new scalable SDF-based renderer and SDF loss, and introducing a new texture refinement network. Comprehensive evaluations and ablations demonstrate the effectiveness of these design choices and state-of-the-art performance.

307

References

- 308 [1] Dejan Azinović, Ricardo Martin-Brualla, Dan B Goldman, Matthias Nießner, and Justus Thies.
309 Neural rgb-d surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer
310 Vision and Pattern Recognition (CVPR)*, pages 6290–6301, June 2022.
- 311 [2] P. Beckmann and A. Spizzichino. *The Scattering of Electromagnetic Waves from Rough
312 Surfaces*. Pergamon Press, 1963.
- 313 [3] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T. Barron, Ce Liu, and Hendrik P.A.
314 Lensch. NeRD: Neural Reflectance Decomposition from Image Collections. In *2021
315 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- 316 [4] Mark Boss, Varun Jampani, Raphael Braun, Ce Liu, Jonathan T. Barron, and Hendrik P. A.
317 Lensch. Neural-PIL: Neural Pre-Integrated Lighting for Reflectance Decomposition. *arXiv
318 preprint*, 2021.
- 319 [5] Brent Burley. Physically-based shading at disney. Technical report, Disney, 2012.
- 320 [6] Ang Cao, Justin Johnson, Andrea Vedaldi, and David Novotny. Lightplane: Highly-scalable
321 components for neural 3d fields. *arXiv*, 2024.
- 322 [7] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello,
323 Orazio Gallo, Leonidas J. Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and
324 Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *Proc.
325 CVPR*, 2022.
- 326 [8] Eric R. Chan, Koki Nagano, Matthew A. Chan, Alexander W. Bergman, Jeong Joon Park, Axel
327 Levy, Miika Aittala, Shalini De Mello, Tero Karras, and Gordon Wetzstein. Generative novel
328 view synthesis with 3D-aware diffusion models. *arXiv.cs*, abs/2304.02602, 2023.
- 329 [9] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. TensoRF: Tensorial radiance
330 fields. In *arXiv*, 2022.
- 331 [10] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3D: Disentangling geometry
332 and appearance for high-quality text-to-3D content creation: Disentangling geometry and
333 appearance for high-quality text-to-3d content creation. *arXiv.cs*, abs/2303.13873, 2023.
- 334 [11] Yabo Chen, Jiemin Fang, Yuyang Huang, Taoran Yi, Xiaopeng Zhang, Lingxi Xie, Xinggang
335 Wang, Wenrui Dai, Hongkai Xiong, and Qi Tian. Cascade-Zero123: One image to highly
336 consistent 3D with self-prompted nearby views. *arXiv.cs*, abs/2312.04424, 2023.
- 337 [12] Zilong Chen, Feng Wang, and Huaping Liu. Text-to-3D using Gaussian splatting. *arXiv*,
338 (2309.16585), 2023.
- 339 [13] Zilong Chen, Yikai Wang, Feng Wang, Zhengyi Wang, and Huaping Liu. V3D: Video diffusion
340 models are effective 3D generators. *arXiv*, 2403.06738, 2024.
- 341 [14] Christopher B. Choy, Danfei Xu, Jun Young Gwak, Kevin Chen, and Silvio Savarese. 3D-R2N2:
342 A unified approach for single and multi-view 3D object reconstruction. In *Proc. ECCV*, 2016.
- 343 [15] Robert L. Cook and Kenneth E. Torrance. A reflectance model for computer graphics. In
344 Doug Green, Tony Lucido, and Henry Fuchs, editors, *Proc. SIGGRAPH*, 1981.
- 345 [16] Xiaoliang Dai, Ji Hou, Chih-Yao Ma, Sam S. Tsai, Jialiang Wang, Rui Wang, Peizhao Zhang,
346 Simon Vandenhende, Xiaofang Wang, Abhimanyu Dubey, Matthew Yu, Abhishek Kadian,
347 Filip Radenovic, Dhruv Mahajan, Kunpeng Li, Yue Zhao, Vladan Petrovic, Mitesh Kumar
348 Singh, Simran Motwani, Yi Wen, Yiwen Song, Roshan Sumbaly, Vignesh Ramanathan, Zijian
349 He, Peter Vajda, and Devi Parikh. Emu: Enhancing image generation models using photogenic
350 needles in a haystack. *CoRR*, abs/2309.15807, 2023.
- 351 [17] François Darmon, Bénédicte Basclle, Jean-Clément Devaux, Pascal Monasse, and Mathieu
352 Aubry. Improving neural implicit surfaces geometry with patch warping. In *Proc. CVPR*,
353 2022.

- 354 [18] Laura Downs, Anthony Francis, Nate Koenig, Brandon Kinman, Ryan Hickman, Krista
 355 Reymann, Thomas B. McHugh, and Vincent Vanhoucke. Google Scanned Objects: A high-
 356 quality dataset of 3D scanned household items. In *Proc. ICRA*, 2022.
- 357 [19] James D. Foley, Andries van Dam, Steven Feiner, and John F. Hughes. *Computer graphics -*
 358 *principles and practice, 3rd Edition*. Addison-Wesley, 2013.
- 359 [20] Qiancheng Fu, Qingshan Xu, Yew-Soon Ong, and Wenbing Tao. Geo-Neus: Geometry-
 360 Consistent Neural Implicit Surfaces Learning for Multi-view Reconstruction. In *NeurIPS*,
 361 2022.
- 362 [21] Jun Gao, Wenzheng Chen, Tommy Xiang, Clement Fuji Tsang, Alec Jacobson, Morgan
 363 McGuire, and Sanja Fidler. Learning deformable tetrahedral meshes for 3D reconstruction. In
 364 *Proc. NeurIPS*, 2020.
- 365 [22] Ruiqi Gao, Aleksander Holynski, Philipp Henzler, Arthur Brussee, Ricardo Martin-Brualla,
 366 Pratul Srinivasan, Jonathan T. Barron, and Ben Poole. CAT3D: Create Anything in 3D with
 367 Multi-View Diffusion Models. *arXiv.cs*, 2024.
- 368 [23] Shubham Goel, Georgia Gkioxari, and Jitendra Malik. Differentiable Stereopsis: Meshes from
 369 multiple views using differentiable rendering. In *CVPR*, 2022.
- 370 [24] Antoine Guédon and Vincent Lepetit. SuGaR: Surface-aligned Gaussian splatting for efficient
 371 3D mesh reconstruction and high-quality mesh rendering. *arXiv.cs*, abs/2311.12775, 2023.
- 372 [25] Anchit Gupta, Wenhan Xiong, Yixin Nie, Ian Jones, and Barlas Oguz. 3DGen: Triplane latent
 373 diffusion for textured mesh generation. *corr*, abs/2303.05371, 2023.
- 374 [26] Jon Hasselgren, Nikolai Hofmann, and Jacob Munkberg. Shape, Light, and Material De-
 375 composition from Images using Monte Carlo Rendering and Denoising. *arXiv preprint*,
 376 2022.
- 377 [27] Lukas Höllerin, Aljaž Božič, Norman Müller, David Novotny, Hung-Yu Tseng, Christian
 378 Richardt, Michael Zollhöfer, and Matthias Nießner. ViewDiff: 3D-Consistent Image Genera-
 379 tion with Text-to-Image Models. *arXiv preprint*, 2024.
- 380 [28] Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan
 381 Sunkavalli, Trung Bui, and Hao Tan. LRM: Large reconstruction model for single image to
 382 3D. In *Proc. ICLR*, 2024.
- 383 [29] Yukun Huang, Jianan Wang, Yukai Shi, Xianbiao Qi, Zheng-Jun Zha, and Lei Zhang.
 384 Dreamtime: An improved optimization strategy for text-to-3D content creation. *CoRR*,
 385 abs/2306.12422, 2023.
- 386 [30] Yifan Jiang, Hao Tang, Jen-Hao Rick Chang, Liangchen Song, Zhangyang Wang, and Lian-
 387 gliang Cao. Efficient-3Dim: Learning a generalizable single-image novel-view synthesizer in
 388 one day. *arXiv*, 2023.
- 389 [31] Yingwenqi Jiang, Jiadong Tu, Yuan Liu, Xifeng Gao, Xiaoxiao Long, Wenping Wang, and
 390 Yuexin Ma. GaussianShader: 3D Gaussian splatting with shading functions for reflective
 391 surfaces. *arXiv.cs*, abs/2311.17977, 2023.
- 392 [32] Heewoo Jun and Alex Nichol. Shape-E: Generating conditional 3D implicit functions. *arXiv*,
 393 2023.
- 394 [33] Philip Torr Junlin Han, Filippos Kokkinos. Vfusion3d: Learning scalable 3d generative models
 395 from video diffusion models. *arXiv preprint*, 2024.
- 396 [34] Angjoo Kanazawa, Shubham Tulsiani, Alexei A. Efros, and Jitendra Malik. Learning category-
 397 specific mesh reconstruction from image collections. In *Proc. ECCV*, 2018.
- 398 [35] Animesh Karnewar, Andrea Vedaldi, David Novotny, and Niloy Mitra. HoloDiffusion: training
 399 a 3D diffusion model using 2D images. In *Proceedings of the IEEE Conference on Computer
 400 Vision and Pattern Recognition (CVPR)*, 2023.

- 401 [36] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3D gaussian
402 splatting for real-time radiance field rendering. *Proc. SIGGRAPH*, 42(4), 2023.
- 403 [37] Seungwook Kim, Yichun Shi, Kejie Li, Minsu Cho, and Peng Wang. Multi-view image
404 prompted multi-view diffusion for improved 3D generation. *arXiv*, 2404.17419, 2024.
- 405 [38] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *Proc. ICLR*,
406 2015.
- 407 [39] Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan
408 Sunkavalli, Greg Shakhnarovich, and Sai Bi. Instant3D: Fast text-to-3D with sparse-view
409 generation and large reconstruction model. *Proc. ICLR*, 2024.
- 410 [40] Zhihao Liang, Qi Zhang, Ying Feng, Ying Shan, and Kui Jia. GS-IR: 3D Gaussian splatting
411 for inverse rendering. *arXiv.cs*, abs/2311.16473, 2023.
- 412 [41] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang,
413 Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3D: High-resolution
414 text-to-3D content creation. *arXiv.cs*, abs/2211.10440, 2022.
- 415 [42] Minghua Liu, Ruoxi Shi, Linghao Chen, Zhuoyang Zhang, Chao Xu, Xinyue Wei, Hansheng
416 Chen, Chong Zeng, Jiayuan Gu, and Hao Su. One-2-3-45++: Fast single image to 3D objects
417 with consistent multi-view generation and 3D diffusion. *arXiv.cs*, abs/2311.07885, 2023.
- 418 [43] Minghua Liu, Chao Xu, Haian Jin, Linghao Chen, Mukund Varma T, Zexiang Xu, and Hao Su.
419 One-2-3-45: Any single image to 3D mesh in 45 seconds without per-shape optimization. In
420 *Proc. NeurIPS*, 2023.
- 421 [44] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl
422 Vondrick. Zero-1-to-3: Zero-shot one image to 3D object. In *Proc. ICCV*, 2023.
- 423 [45] Shichen Liu, Tianye Li, Weikai Chen, and Hao Li. Soft rasterizer: A differentiable renderer
424 for image-based 3D reasoning. *arXiv.cs*, abs/1904.01786, 2019.
- 425 [46] Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping
426 Wang. SyncDreamer: Generating multiview-consistent images from a single-view image.
427 *arXiv*, (2309.03453), 2023.
- 428 [47] Zexiang Liu, Yangguang Li, Youtian Lin, Xin Yu, Sida Peng, Yan-Pei Cao, Xiaojuan Qi,
429 Xiaoshui Huang, Ding Liang, and Wanli Ouyang. UniDream: Unifying Diffusion Priors for
430 Relightable Text-to-3D Generation. *arXiv preprint*, 2023.
- 431 [48] Xiaoxiao Long, Yuanchen Guo, Cheng Lin, Yuan Liu, Zhiyang Dou, Lingjie Liu, Yuexin
432 Ma, Song-Hai Zhang, Marc Habermann, Christian Theobalt, and Wenping Wang. Wonder3D:
433 Single image to 3D using cross-domain diffusion. *arXiv.cs*, abs/2310.15008, 2023.
- 434 [49] W. Lorensen and H. Cline. Marching cubes: A high resolution 3D surface construction
435 algorithm. *ACM Computer Graphocs*, 21(24), 1987.
- 436 [50] Jonathan Lorraine, Kevin Xie, Xiaohui Zeng, Chen-Hsuan Lin, Towaki Takikawa, Nicholas
437 Sharp, Tsung-Yi Lin, Ming-Yu Liu, Sanja Fidler, and James Lucas. ATT3D: amortized
438 text-to-3D object synthesis. In *Proc. ICCV*, 2023.
- 439 [51] Tiange Luo, Chris Rockwell, Honglak Lee, and Justin Johnson. Scalable 3d captioning with
440 pretrained models. *arXiv preprint*, 2023.
- 441 [52] Luke Melas-Kyriazi, Iro Laina, Christian Rupprecht, Natalia Neverova, Andrea Vedaldi, Oran
442 Gafni, and Filippos Kokkinos. IM-3D: Iterative multiview diffusion and reconstruction for
443 high-quality 3D generation. *arXiv preprint*, (abs/2402.08682), 2024.
- 444 [53] Luke Melas-Kyriazi, Christian Rupprecht, Iro Laina, and Andrea Vedaldi. RealFusion: 360
445 reconstruction of any object from a single image. In *Proceedings of the IEEE Conference on
446 Computer Vision and Pattern Recognition (CVPR)*, 2023.

- 447 [54] Antoine Mercier, Ramin Nakhli, Mahesh Reddy, and Rajeev Yasarla. HexaGen3D: Stabledif-
448 fusion is just one step away from fast and diverse text-to-3D generation. *arXiv*, 2024.
- 449 [55] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas
450 Geiger. Occupancy Networks: Learning 3D Reconstruction in Function Space. In *CVPR*,
451 2019.
- 452 [56] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi,
453 and Ren Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. In *Proc.
454 ECCV*, 2020.
- 455 [57] Tom Monnier, Jake Austin, Angjoo Kanazawa, Alexei A. Efros, and Mathieu Aubry. Differentiable
456 blocks world: Qualitative 3d decomposition by rendering primitives. *abs/2307.05473*,
457 2023.
- 458 [58] Tom Monnier, Matthew Fisher, Alexei A. Efros, and Mathieu Aubry. Share With Thy Neigh-
459 bors: Single-View Reconstruction by Cross-Instance Consistency. In *ECCV*, 2022.
- 460 [59] Norman Müller, Yawar Siddiqui, Lorenzo Porzi, Samuel Rota Bulò, Peter Kontschieder, and
461 Matthias Nießner. DiffRF: Rendering-guided 3D radiance field diffusion. In *Proc. CVPR*,
462 2023.
- 463 [60] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics
464 primitives with a multiresolution hash encoding. In *Proc. SIGGRAPH*, 2022.
- 465 [61] Jacob Munkberg, Wenzheng Chen, Jon Hasselgren, Alex Evans, Tianchang Shen, Thomas
466 Muller, Jun Gao, and Sanja Fidler. Extracting Triangular 3D Models, Materials, and Lighting
467 From Images. In *CVPR*, 2022.
- 468 [62] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. Point-E:
469 A system for generating 3D point clouds from complex prompts. *arXiv.cs*, *abs/2212.08751*,
470 2022.
- 471 [63] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable
472 Volumetric Rendering: Learning Implicit 3D Representations without 3D Supervision. In
473 *CVPR*, 2020.
- 474 [64] Michael Oechsle, Songyou Peng, and Andreas Geiger. UNISURF: unifying neural implicit
475 surfaces and radiance fields for multi-view reconstruction. *arXiv.cs*, *abs/2104.10078*, 2021.
- 476 [65] OpenAI. Triton: Open-source gpu programming for neural networks. <https://github.com/triton-lang/triton>.
- 477 [66] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. DreamFusion: Text-to-3D
478 using 2D diffusion. In *Proc. ICLR*, 2023.
- 479 [67] Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren, Aliaksandr Siarohin, Bing Li, Hsin-
480 Ying Lee, Ivan Skorokhodov, Peter Wonka, Sergey Tulyakov, and Bernard Ghanem. Magic123:
481 One image to high-quality 3D object generation using both 2D and 3D diffusion priors.
482 *arXiv.cs*, *abs/2306.17843*, 2023.
- 483 [68] Lingteng Qiu, Guanying Chen, Xiaodong Gu, Qi Zuo, Mutian Xu, Yushuang Wu, Weihao Yuan,
484 Zilong Dong, Liefeng Bo, and Xiaoguang Han. Richdreamer: A generalizable normal-depth
485 diffusion model for detail richness in text-to-3D. *arXiv.cs*, *abs/2311.16918*, 2023.
- 486 [69] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin
487 Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv*, 2020.
- 488 [70] Christophe Schlick. An inexpensive BRDF model for physically-based rendering. *Comput.
489 Graph. Forum*, 13(3), 1994.
- 490 [71] Ruoxi Shi, Hansheng Chen, Zhuoyang Zhang, Minghua Liu, Chao Xu, Xinyue Wei, Linghao
491 Chen, Chong Zeng, and Hao Su. Zero123++: a single image to consistent multi-view diffusion
492 base model. *arXiv.cs*, *abs/2310.15110*, 2023.

- 494 [72] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. MVDream:
 495 Multi-view diffusion for 3D generation. In *Proc. ICLR*, 2024.
- 496 [73] J. Ryan Shue, Eric Ryan Chan, Ryan Po, Zachary Ankner, Jiajun Wu, and Gordon Wetzstein.
 497 3D neural field generation using triplane diffusion. *arXiv.cs*, abs/2211.16677, 2022.
- 498 [74] B. Smith. Geometrical shadowing of a random rough surface. *IEEE Trans. on Antennas and
 499 Propagation*, 15(5), 1967.
- 500 [75] Jingxiang Sun, Bo Zhang, Ruizhi Shao, Lizhen Wang, Wen Liu, Zhenda Xie, and Yebin
 501 Liu. DreamCraft3D: Hierarchical 3D generation with bootstrapped diffusion prior. *arXiv.cs*,
 502 abs/2310.16818, 2023.
- 503 [76] Stanislaw Szymanowicz, Christian Rupprecht, and Andrea Vedaldi. Viewset diffusion: (0-
 504)image-conditioned 3D generative models from 2D data. In *Proceedings of the International
 505 Conference on Computer Vision (ICCV)*, 2023.
- 506 [77] Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, Gang Zeng, and Ziwei Liu. LGM:
 507 Large multi-view Gaussian model for high-resolution 3D content creation. *arXiv*, 2402.05054,
 508 2024.
- 509 [78] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. DreamGaussian: Genera-
 510 tive gaussian splatting for efficient 3D content creation. *arXiv*, (2309.16653), 2023.
- 511 [79] Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi, Lizhuang Ma, and Dong Chen.
 512 Make-It-3D: High-fidelity 3d creation from A single image with diffusion prior. *arXiv.cs*,
 513 abs/2303.14184, 2023.
- 514 [80] Shitao Tang, Jiacheng Chen, Dilin Wang, Chengzhou Tang, Fuyang Zhang, Yuchen Fan,
 515 Vikas Chandra, Yasutaka Furukawa, and Rakesh Ranjan. MVDiffusion++: A dense high-
 516 resolution multi-view diffusion model for single or sparse-view 3d object reconstruction. *arXiv*,
 517 2402.12712, 2024.
- 518 [81] Luma Team. Luma genie 1.0. <https://www.luma-ai.com/luma-genie-1-0/>.
- 519 [82] Meshy Team. Meshy - AI 3D Model Generator with pbr materials— meshy.ai. <https://www.meshy.ai/>.
- 520 [83] Dmitry Tochilkin, David Pankratz, Zexiang Liu, Zixuan Huang, Adam Letts, Yangguang Li,
 521 Ding Liang, Christian Laforte, Varun Jampani, and Yan-Pei Cao. TripoSR: fast 3D object
 522 reconstruction from a single image. 2403.02151, 2024.
- 523 [84] K. E. Torrance and E. M. Sparrow. Theory for off-specular reflection from roughened surfaces.
 524 *J. Opt. Soc. Am.*, 57(9), 1967.
- 525 [85] Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine
 526 Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, D. Bikel,
 527 Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude
 528 Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman
 529 Goyal, A. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez,
 530 Madian Khabsa, Isabel M. Kloumann, A. Korenev, Punit Singh Koura, Marie-Anne Lachaux,
 531 Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor
 532 Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein,
 533 Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian,
 534 Xia Tan, Bin Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu
 535 Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien
 536 Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation
 537 and fine-tuned chat models, 7 2023.
- 538 [86] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
 539 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017.
- 540 [87] Bruce Walter, Stephen R. Marschner, Hongsong Li, and Kenneth E. Torrance. Microfacet
 541 models for refraction through rough surfaces. In *Proc. Eurographics*, 2007.

- 543 [88] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A. Yeh, and Greg Shakhnarovich. Score
 544 Jacobian Chaining: Lifting Pretrained 2D Diffusion Models for 3D Generation. In *CVPR*,
 545 2023.
- 546 [89] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang.
 547 NeuS: Learning neural implicit surfaces by volume rendering for multi-view reconstruction.
 548 *arXiv.cs*, abs/2106.10689, 2021.
- 549 [90] Peng Wang and Yichun Shi. ImageDream: Image-prompt multi-view diffusion for 3D genera-
 550 tion. In *Proc. ICLR*, 2024.
- 551 [91] Tengfei Wang, Bo Zhang, Ting Zhang, Shuyang Gu, Jianmin Bao, Tadas Baltrusaitis, Jingjing
 552 Shen, Dong Chen, Fang Wen, Qifeng Chen, and Baining Guo. Rodin: A generative model for
 553 sculpting 3D digital avatars using diffusion. In *Proc. CVPR*, 2023.
- 554 [92] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Pro-
 555 lificDreamer: High-fidelity and diverse text-to-3D generation with variational score distillation.
 556 *arXiv.cs*, abs/2305.16213, 2023.
- 557 [93] Zhengyi Wang, Yikai Wang, Yifei Chen, Chendong Xiang, Shuo Chen, Dajiang Yu, Chongxuan
 558 Li, Hang Su, and Jun Zhu. CRM: Single image to 3D textured mesh with convolutional
 559 reconstruction model. *arXiv*, (2403.05034), 2024.
- 560 [94] Xinyue Wei, Kai Zhang, Sai Bi, Hao Tan, Fujun Luan, Valentin Deschaintre, Kalyan Sunkavalli,
 561 Hao Su, and Zexiang Xu. MeshLRM: large reconstruction model for high-quality mesh. *arXiv*,
 562 2404.12385, 2024.
- 563 [95] Haohan Weng, Tianyu Yang, Jianan Wang, Yu Li, Tong Zhang, C. L. Philip Chen, and Lei
 564 Zhang. Consistent123: Improve consistency for one image to 3D object synthesis. *arXiv*,
 565 2023.
- 566 [96] Rundi Wu, Ben Mildenhall, Philipp Henzler, Keunhong Park, Ruiqi Gao, Daniel Watson,
 567 Pratul P. Srinivasan, Dor Verbin, Jonathan T. Barron, Ben Poole, and Aleksander Holynski.
 568 ReconFusion: 3D Reconstruction with Diffusion Priors. *arXiv preprint*, 2023.
- 569 [97] Shangzhe Wu, Christian Rupprecht, and Andrea Vedaldi. Unsupervised learning of probably
 570 symmetric deformable 3D objects from images in the wild. In *Proceedings of the IEEE*
 571 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- 572 [98] Kevin Xie, Jonathan Lorraine, Tianshi Cao, Jun Gao, James Lucas, Antonio Torralba, Sanja
 573 Fidler, and Xiaohui Zeng. LATTE3D: Large-scale amortized text-to-enhanced3D synthesis.
 574 In *arXiv*, 2024.
- 575 [99] Zhang Xiuming, Srinivasan Pratul P., Deng Boyang, Debevec Paul, Freeman William T., and
 576 Barron Jonathan T. NeRFactor: neural factorization of shape and reflectance under an unknown
 577 illumination. In *Proc. SIGGRAPH*, 2021.
- 578 [100] Jiale Xu, Weihao Cheng, Yiming Gao, Xintao Wang, Shenghua Gao, and Ying Shan. In-
 579 stantMesh: efficient 3D mesh generation from a single image with sparse-view large recon-
 580 struction models. *arXiv*, 2404.07191, 2024.
- 581 [101] Xudong Xu, Zhaoyang Lyu, Xingang Pan, and Bo Dai. MATLAKER: Material-Aware Text-to-
 582 3D via LAtent BRDF auto-EncodeR. *arXiv preprint*, 2023.
- 583 [102] Yinghao Xu, Zifan Shi, Wang Yifan, Hansheng Chen, Ceyuan Yang, Sida Peng, Yujun
 584 Shen, and Gordon Wetzstein. GRM: Large gaussian reconstruction model for efficient 3D
 585 reconstruction and generation. *arXiv*, 2403.14621, 2024.
- 586 [103] Yinghao Xu, Hao Tan, Fujun Luan, Sai Bi, Peng Wang, Jiahao Li, Zifan Shi, Kalyan Sunkavalli,
 587 Gordon Wetzstein, Zexiang Xu, and Kai Zhang. DMV3D: Denoising multi-view diffusion
 588 using 3D large reconstruction model. In *Proc. ICLR*, 2024.
- 589 [104] Jiayu Yang, Ziang Cheng, Yunfei Duan, Pan Ji, and Hongdong Li. ConsistNet: Enforcing 3D
 590 consistency for multi-view images diffusion. *arXiv.cs*, abs/2310.10343, 2023.

- 591 [105] Yunhan Yang, Yukun Huang, Xiaoyang Wu, Yuan-Chen Guo, Song-Hai Zhang, Hengshuang
 592 Zhao, Tong He, and Xihui Liu. DreamComposer: Controllable 3D object generation via
 593 multi-view conditions. *arXiv.cs*, abs/2312.03611, 2023.
- 594 [106] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit
 595 surfaces. *arXiv.cs*, abs/2106.12052, 2021.
- 596 [107] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Ronen Basri, and Yaron
 597 Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance.
 598 In *Proc. NeurIPS*, 2020.
- 599 [108] Lior Yariv, Omri Puny, Natalia Neverova, Oran Gafni, and Yaron Lipman. Mosaic-SDF for
 600 3D generative models. *arXiv.cs*, abs/2312.09222, 2023.
- 601 [109] Taoran Yi, Jiemin Fang, Guanjun Wu, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Qi Tian, and
 602 Xinggang Wang. GaussianDreamer: Fast generation from text to 3D gaussian splatting with
 603 point cloud priors. *arXiv.cs*, abs/2310.08529, 2023.
- 604 [110] Jonathan Young. Xatlas: Mesh parameterization / uv unwrapping library, 2022. GitHub
 605 repository.
- 606 [111] Wangbo Yu, Li Yuan, Yan-Pei Cao, Xiangjun Gao, Xiaoyu Li, Long Quan, Ying Shan, and
 607 Yonghong Tian. HiFi-123: Towards high-fidelity one image to 3D content generation. *arXiv.cs*,
 608 abs/2310.06744, 2023.
- 609 [112] Jason Y. Zhang, Gengshan Yang, Shubham Tulsiani, and Deva Ramanan. NeRS: Neural
 610 Reflectance Surfaces for Sparse-view 3D Reconstruction in the Wild. In *NeurIPS*, 2021.
- 611 [113] Kai Zhang, Sai Bi, Hao Tan, Yuanbo Xiangli, Nanxuan Zhao, Kalyan Sunkavalli, and Zexiang
 612 Xu. GS-LRM: large reconstruction model for 3D Gaussian splatting. *arXiv*, 2404.19702, 2024.
- 613 [114] Kai Zhang, Fujun Luan, Qianqian Wang, Kavita Bala, and Noah Snavely. PhySG: Inverse
 614 Rendering with Spherical Gaussians for Physics-based Material Editing and Relighting. *arXiv
 615 preprint*, 2021.
- 616 [115] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unrea-
 617 sonable effectiveness of deep features as a perceptual metric. In *Proc. CVPR*, pages 586–595,
 618 2018.
- 619 [116] Xiaoyu Zhou, Xingjian Ran, Yajiao Xiong, Jinlin He, Zhiwei Lin, Yongtao Wang, Deqing
 620 Sun, and Ming-Hsuan Yang. GALA3D: Towards text-to-3D complex scene generation via
 621 layout-guided generative gaussian splatting. *arXiv.cs*, abs/2402.07207, 2024.
- 622 [117] Junzhe Zhu and Peiye Zhuang. HiFA: High-fidelity text-to-3D with advanced diffusion
 623 guidance. *CoRR*, abs/2305.18766, 2023.
- 624 [118] Zi-Xin Zou, Zhipeng Yu, Yuan-Chen Guo, Yangguang Li, Ding Liang, Yan-Pei Cao, and
 625 Song-Hai Zhang. Triplane meets Gaussian splatting: Fast and generalizable single-view 3D
 626 reconstruction with transformers. *arXiv.cs*, abs/2312.09147, 2023.

627 **NeurIPS Paper Checklist**

628 **1. Claims**

629 Question: Do the main claims made in the abstract and introduction accurately reflect the
630 paper's contributions and scope?

631 Answer: [Yes]

632 Justification: See Sec. 1 and the abstract.

633 Guidelines:

- 634 • The answer NA means that the abstract and introduction do not include the claims
635 made in the paper.
- 636 • The abstract and/or introduction should clearly state the claims made, including the
637 contributions made in the paper and important assumptions and limitations. A No or
638 NA answer to this question will not be perceived well by the reviewers.
- 639 • The claims made should match theoretical and experimental results, and reflect how
640 much the results can be expected to generalize to other settings.
- 641 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
642 are not attained by the paper.

643 **2. Limitations**

644 Question: Does the paper discuss the limitations of the work performed by the authors?

645 Answer: [Yes]

646 Justification: See App. A.2

647 Guidelines:

- 648 • The answer NA means that the paper has no limitation while the answer No means that
649 the paper has limitations, but those are not discussed in the paper.
- 650 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 651 • The paper should point out any strong assumptions and how robust the results are to
652 violations of these assumptions (e.g., independence assumptions, noiseless settings,
653 model well-specification, asymptotic approximations only holding locally). The authors
654 should reflect on how these assumptions might be violated in practice and what the
655 implications would be.
- 656 • The authors should reflect on the scope of the claims made, e.g., if the approach was
657 only tested on a few datasets or with a few runs. In general, empirical results often
658 depend on implicit assumptions, which should be articulated.
- 659 • The authors should reflect on the factors that influence the performance of the approach.
660 For example, a facial recognition algorithm may perform poorly when image resolution
661 is low or images are taken in low lighting. Or a speech-to-text system might not be
662 used reliably to provide closed captions for online lectures because it fails to handle
663 technical jargon.
- 664 • The authors should discuss the computational efficiency of the proposed algorithms
665 and how they scale with dataset size.
- 666 • If applicable, the authors should discuss possible limitations of their approach to
667 address problems of privacy and fairness.
- 668 • While the authors might fear that complete honesty about limitations might be used by
669 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
670 limitations that aren't acknowledged in the paper. The authors should use their best
671 judgment and recognize that individual actions in favor of transparency play an impor-
672 tant role in developing norms that preserve the integrity of the community. Reviewers
673 will be specifically instructed to not penalize honesty concerning limitations.

674 **3. Theory Assumptions and Proofs**

675 Question: For each theoretical result, does the paper provide the full set of assumptions and
676 a complete (and correct) proof?

677 Answer: [NA]

678 Justification: No theoretical results present.

679 Guidelines:

- 680 • The answer NA means that the paper does not include theoretical results.
- 681 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 682 • referenced.
- 683 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 684 • The proofs can either appear in the main paper or the supplemental material, but if
- 685 • they appear in the supplemental material, the authors are encouraged to provide a short
- 686 • proof sketch to provide intuition.
- 687 • Inversely, any informal proof provided in the core of the paper should be complemented
- 688 • by formal proofs provided in appendix or supplemental material.
- 689 • Theorems and Lemmas that the proof relies upon should be properly referenced.

690 **4. Experimental Result Reproducibility**

691 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
692 perimental results of the paper to the extent that it affects the main claims and/or conclusions
693 of the paper (regardless of whether the code and data are provided or not)?

694 Answer: [Yes]

695 Justification: Please see Secs. 3 and 4 and Apps. A.6.2 and A.7 for all technical details
696 required for reproducing the paper results.

697 Guidelines:

- 698 • The answer NA means that the paper does not include experiments.
- 699 • If the paper includes experiments, a No answer to this question will not be perceived
- 700 • well by the reviewers: Making the paper reproducible is important, regardless of
- 701 • whether the code and data are provided or not.
- 702 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 703 • to make their results reproducible or verifiable.
- 704 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 705 • For example, if the contribution is a novel architecture, describing the architecture fully
- 706 • might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 707 • be necessary to either make it possible for others to replicate the model with the same
- 708 • dataset, or provide access to the model. In general, releasing code and data is often
- 709 • one good way to accomplish this, but reproducibility can also be provided via detailed
- 710 • instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 711 • of a large language model), releasing of a model checkpoint, or other means that are
- 712 • appropriate to the research performed.
- 713 • While NeurIPS does not require releasing code, the conference does require all submis-
- 714 • sions to provide some reasonable avenue for reproducibility, which may depend on the
- 715 • nature of the contribution. For example
 - 716 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
 - 717 • to reproduce that algorithm.
 - 718 (b) If the contribution is primarily a new model architecture, the paper should describe
 - 719 • the architecture clearly and fully.
 - 720 (c) If the contribution is a new model (e.g., a large language model), then there should
 - 721 • either be a way to access this model for reproducing the results or a way to reproduce
 - 722 • the model (e.g., with an open-source dataset or instructions for how to construct
 - 723 • the dataset).
 - 724 (d) We recognize that reproducibility may be tricky in some cases, in which case
 - 725 • authors are welcome to describe the particular way they provide for reproducibility.
 - 726 • In the case of closed-source models, it may be that access to the model is limited in
 - 727 • some way (e.g., to registered users), but it should be possible for other researchers
 - 728 • to have some path to reproducing or verifying the results.

729 **5. Open access to data and code**

730 Question: Does the paper provide open access to the data and code, with sufficient instruc-

731 tions to faithfully reproduce the main experimental results, as described in supplemental

732 material?

733 Answer: [No]

734 Justification: We do not currently plan to open-source the method or training data.

735 Guidelines:

- 736 • The answer NA means that paper does not include experiments requiring code.
- 737 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 738 • While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- 739 • The instructions should contain the exact command and environment needed to run to 740 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 741 • The authors should provide instructions on data access and preparation, including how 742 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 743 • The authors should provide scripts to reproduce all experimental results for the new 744 proposed method and baselines. If only a subset of experiments are reproducible, they 745 should state which ones are omitted from the script and why.
- 746 • At submission time, to preserve anonymity, the authors should release anonymized 747 versions (if applicable).
- 748 • Providing as much information as possible in supplemental material (appended to the 749 paper) is recommended, but including URLs to data and code is permitted.
- 750 • Providing as much information as possible in supplemental material (appended to the 751 paper) is recommended, but including URLs to data and code is permitted.

755 6. Experimental Setting/Details

756 Question: Does the paper specify all the training and test details (e.g., data splits, hyper- 757 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the 758 results?

759 Answer: [Yes]

760 Justification: Please see Secs. 3 and 4 and Apps. A.6.2 and A.7 for all technical details 761 required for reproducing the paper results.

762 Guidelines:

- 763 • The answer NA means that the paper does not include experiments.
- 764 • The experimental setting should be presented in the core of the paper to a level of detail 765 that is necessary to appreciate the results and make sense of them.
- 766 • The full details can be provided either with the code, in appendix, or as supplemental 767 material.

768 7. Experiment Statistical Significance

769 Question: Does the paper report error bars suitably and correctly defined or other appropriate 770 information about the statistical significance of the experiments?

771 Answer: [No]

772 Justification: Reported metrics have negligible or zero variance.

773 Guidelines:

- 774 • The answer NA means that the paper does not include experiments.
- 775 • The authors should answer "Yes" if the results are accompanied by error bars, confi- 776 dence intervals, or statistical significance tests, at least for the experiments that support 777 the main claims of the paper.
- 778 • The factors of variability that the error bars are capturing should be clearly stated (for 779 example, train/test split, initialization, random drawing of some parameter, or overall 780 run with given experimental conditions).
- 781 • The method for calculating the error bars should be explained (closed form formula, 782 call to a library function, bootstrap, etc.)
- 783 • The assumptions made should be given (e.g., Normally distributed errors).

- 784 • It should be clear whether the error bar is the standard deviation or the standard error
 785 of the mean.
 786 • It is OK to report 1-sigma error bars, but one should state it. The authors should
 787 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
 788 of Normality of errors is not verified.
 789 • For asymmetric distributions, the authors should be careful not to show in tables or
 790 figures symmetric error bars that would yield results that are out of range (e.g. negative
 791 error rates).
 792 • If error bars are reported in tables or plots, The authors should explain in the text how
 793 they were calculated and reference the corresponding figures or tables in the text.

794 **8. Experiments Compute Resources**

795 Question: For each experiment, does the paper provide sufficient information on the com-
 796 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 797 the experiments?

798 Answer: [Yes]

799 Justification: See Sec. 4 and App. A.6.2

800 Guidelines:

- 801 • The answer NA means that the paper does not include experiments.
 802 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
 803 or cloud provider, including relevant memory and storage.
 804 • The paper should provide the amount of compute required for each of the individual
 805 experimental runs as well as estimate the total compute.
 806 • The paper should disclose whether the full research project required more compute
 807 than the experiments reported in the paper (e.g., preliminary or failed experiments that
 808 didn't make it into the paper).

809 **9. Code Of Ethics**

810 Question: Does the research conducted in the paper conform, in every respect, with the
 811 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

812 Answer: [Yes]

813 Justification: [NA]

814 Guidelines:

- 815 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 816 • If the authors answer No, they should explain the special circumstances that require a
 817 deviation from the Code of Ethics.
 818 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
 819 eration due to laws or regulations in their jurisdiction).

820 **10. Broader Impacts**

821 Question: Does the paper discuss both potential positive societal impacts and negative
 822 societal impacts of the work performed?

823 Answer: [Yes],

824 Justification: See App. A.1

825 Guidelines:

- 826 • The answer NA means that there is no societal impact of the work performed.
 827 • If the authors answer NA or No, they should explain why their work has no societal
 828 impact or why the paper does not address societal impact.
 829 • Examples of negative societal impacts include potential malicious or unintended uses
 830 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
 831 (e.g., deployment of technologies that could make decisions that unfairly impact specific
 832 groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Data or models are not released

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We use an internal dataset of artist-created meshes whose use has been approved by a professional legal team.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- 886 • If this information is not available online, the authors are encouraged to reach out to
887 the asset's creators.

888 **13. New Assets**

889 Question: Are new assets introduced in the paper well documented and is the documentation
890 provided alongside the assets?

891 Answer: [NA]

892 Justification:

893 Guidelines:

- 894 • The answer NA means that the paper does not release new assets.
- 895 • Researchers should communicate the details of the dataset/code/model as part of their
896 submissions via structured templates. This includes details about training, license,
897 limitations, etc.
- 898 • The paper should discuss whether and how consent was obtained from people whose
899 asset is used.
- 900 • At submission time, remember to anonymize your assets (if applicable). You can either
901 create an anonymized URL or include an anonymized zip file.

902 **14. Crowdsourcing and Research with Human Subjects**

903 Question: For crowdsourcing experiments and research with human subjects, does the paper
904 include the full text of instructions given to participants and screenshots, if applicable, as
905 well as details about compensation (if any)?

906 Answer: [Yes]

907 Justification: See App. A.4.

908 Guidelines:

- 909 • The answer NA means that the paper does not involve crowdsourcing nor research with
910 human subjects.
- 911 • Including this information in the supplemental material is fine, but if the main contribu-
912 tion of the paper involves human subjects, then as much detail as possible should be
913 included in the main paper.
- 914 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
915 or other labor should be paid at least the minimum wage in the country of the data
916 collector.

917 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
918 Subjects**

919 Question: Does the paper describe potential risks incurred by study participants, whether
920 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
921 approvals (or an equivalent approval/review based on the requirements of your country or
922 institution) were obtained?

923 Answer: [Yes]

924 Justification: The internal legal and privacy audit has approved our user study.

925 Guidelines:

- 926 • The answer NA means that the paper does not involve crowdsourcing nor research with
927 human subjects.
- 928 • Depending on the country in which research is conducted, IRB approval (or equivalent)
929 may be required for any human subjects research. If you obtained IRB approval, you
930 should clearly state this in the paper.
- 931 • We recognize that the procedures for this may vary significantly between institutions
932 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
933 guidelines for their institution.
- 934 • For initial submissions, do not include any information that would break anonymity (if
935 applicable), such as the institution conducting the review.

936 **A Appendix**

937 **A.1 Societal Impact**

938 Generative models pose a societal threat, and we do not condone the use of our work for generating
939 deep fakes, spreading misinformation, or creating inappropriate content. Safeguards should be
940 implemented to prevent abuse, such as filtering input text prompts and detecting unsafe content in
941 generated 3D models. Additionally, our generation process may be vulnerable to biases present in
942 the data and 3D models it relies on, potentially perpetuating these biases in the generated content.
943 Despite these risks, our method can augment the work of artists and creative professionals by serving
944 as a complementary tool to boost productivity. It also holds the potential to democratize 3D content
945 creation, making it accessible to those without specialized knowledge or expensive proprietary
946 software.

947 **A.2 Limitations**

948 Emu3D significantly advances shape generation but faces several limitations. Despite the fine-tuning
949 of the multiview image grid generator for view consistency, it is not guaranteed, potentially impacting
950 3D reconstruction quality. Since we use an SDF as an underlying representation, the reconstructor
951 may incorrectly model translucent objects or thin structures like hair or fur. Additionally, while our
952 scalable Triton [65] implementation supports a triplane representation at a resolution of 128×128 ,
953 this representation is inefficient, as much of its capacity is used for empty regions. Future work
954 could explore scalable representations such as octrees, sparse voxel grids, and hash-based methods,
955 which may remove the need for a separate texture enhancement model. We also only predict albedo,
956 metallicity and roughness, and not emissivity or ambient occlusions. Finally, our method has only
957 been tested on object-level reconstructions, leaving scene-scale 3D generation for future research.

958 **A.3 Additional qualitative comparisons**

959 This section describes additional qualitative comparisons that, due to limited space, could not be
960 included in the main paper. Firstly, please refer to the video attached in the supplementary material
961 which provides a holistic presentation of Emu3D’s qualitative results. In Fig. 6, we highlight the
962 contributions of EmuLRM in geometry, texture and material reconstruction. In Fig. 11, we visualize
963 the control of materials provided by Emu3D, i.e., metallicity and roughness, by changing the text
964 prompt for the same concept. Fig. 7 visualizes the renders of the material maps extracted with
965 EmuLRM given four input test views. In Fig. 8, we provide a more extensive qualitative comparison
966 to MeshLRM, the strongest few-view reconstruction baseline. Finally, Fig. 9 provides a gallery of
967 text-conditioned generations depicting Blender-shaded renders together with the rendered PBR maps.
968

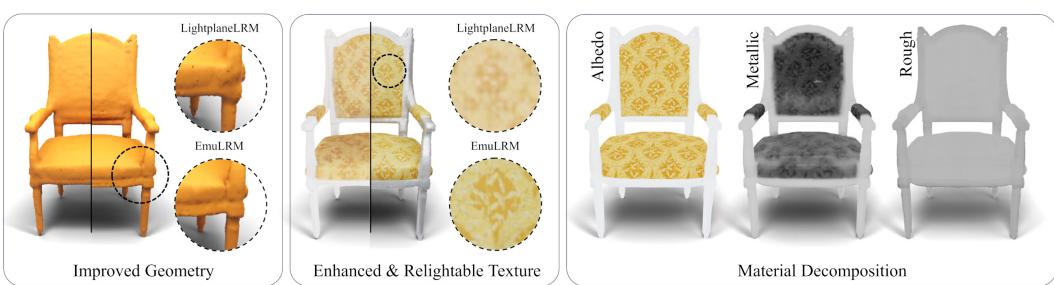


Figure 6: EmuLRM builds upon LightplaneLRM [6], providing improved geometry by employing SDF as a representation, along with direct scalable losses in 3D, improved texture using a UV space texture refiner, and material decomposition by predicting material properties regularized through a novel deferred shading loss.

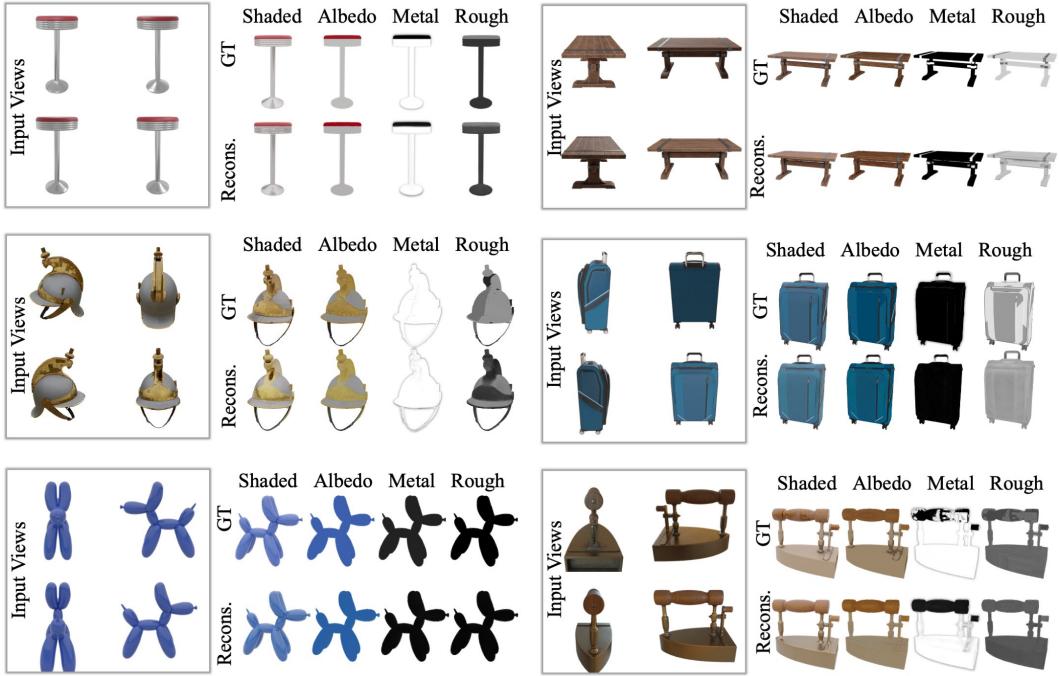


Figure 7: Sparse view reconstruction with intrinsic decomposition. Here the EmuLRM takes 4 shaded views as input and reconstructs the 3D object along with its albedo, metallic and roughness properties.

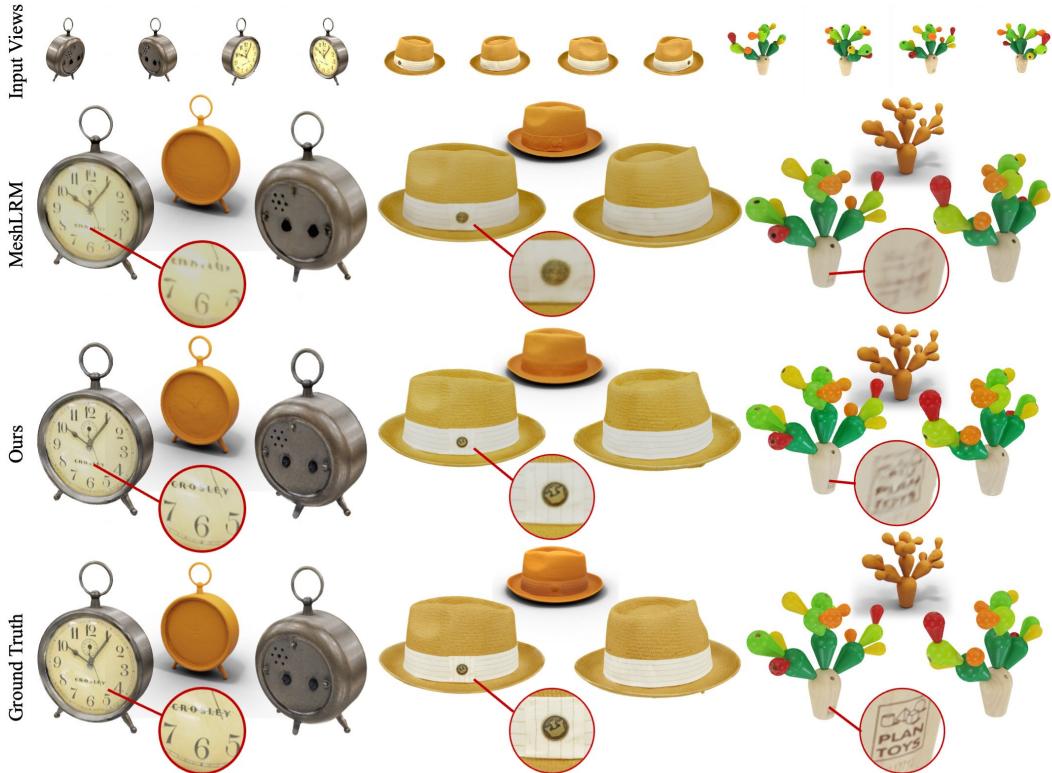


Figure 8: Qualitative comparison on the task of sparse view reconstruction against MeshLRM [94]. Note the higher quality texture detail in Emu3D, as seen in the text on the clock, the logo on the hat, and the text on the plant base. Since an open-source implementation of MeshLRM has not yet been released, we compare against the meshes provided on their webpage.



Figure 9: **Text-to-3D** meshes generated by Emu3D along with their PBR decomposition. Note that Emu3D provides detailed albedo and material properties, as highlighted by the metallicity of the platter (top right) and the golden objects (last row).

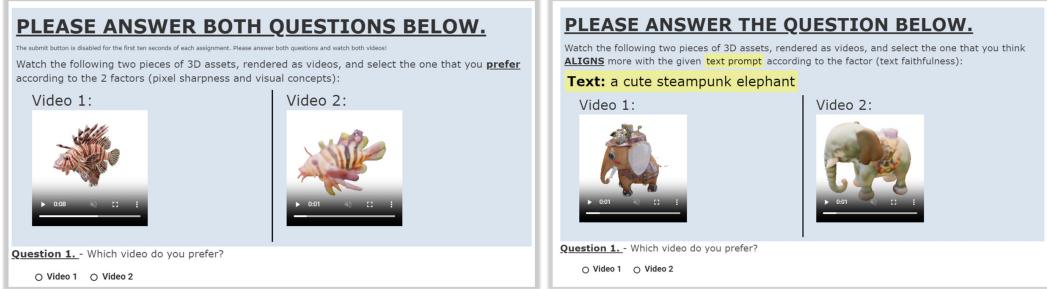


Figure 10: User study interface submitted to a standard crowdsourcing marketplace. Participants are shown videos corresponding to Emu3D and a baseline in a random order, and asked their preference in terms of either quality (left) or faithfulness to the text prompt (right).

969 A.4 User-study details

970 As described in Sec. 4.2, we conducted an user study on 404 meshes generated using the DreamFu-
 971 sion [66] prompt-set on a standard crowdsourcing marketplace. In the study, users were shown 360°
 972 videos of the generated and baseline meshes and were asked to rate them based on 3D shape quality
 973 and alignment with the text prompt as shown in Fig. 10. They were asked to consider various factors
 974 like identity (whether the object matches what is described in the prompt), texture, existence of Janus
 975 problems, and bad geometry (like floaters, disconnected components, etc). 11,080 responses were
 976 collected in total, with 5 responses per pair of videos, to eliminate variance in user preference.

977 A.5 Additional text-to-3d comparisons

978 While Tab. 2 compared Emu3D’s text-to-3d gen-
 979 erations to several fast baselines, for comple-
 980 teness, this section includes additional compari-
 981 sons to significantly slower methods. More
 982 specifically, we conduct the same user study as
 983 in Tab. 2 but we compare to the “refinement”
 984 stages of the industry baselines Meshy v3 and
 985 Luma Genie whose asset generation time is 5
 986 and 10 min respectively. Tab. 4 contains the
 987 results of our user-study. Emu3D significantly out-
 988 performs Meshy in both text fidelity and visual
 989 quality while being 10× faster. Surprisingly,
 990 Emu3D is on par with Luma Genie in text fidelity and wins in 40% of cases in visual quality. This is
 991 a remarkable result considering Emu3D’s 20× better generation time.

992 A.6 Additional technical details

993 A.6.1 Grid Generator

994 We employ a text-to-image diffusion model pre-trained on billions of images annotated with text [16]
 995 and expand its input and output channels by a factor of 2 to support simultaneous generation of shaded
 996 appearance and albedo. We finetune the model to predict a grid of four images I_i , $i = 1, \dots, 4$, in
 997 similar fashion to [72, 39] via minimization of the standard diffusion loss. Training spans a total of 2
 998 days, employing 32 A100 GPUs with a total batch size of 128 and a learning rate of 10^{-5} .

999 A.6.2 EmuLRM

As mentioned in the main paper, EmuLRM is optimized using the direct SDF loss \mathcal{L}_{sdf} , PBR loss \mathcal{L}_{pbr} , deferred shading loss \mathcal{L}_{def} , the binary cross-entropy mask loss $\mathcal{L}_{\text{mask}}$, and the depth-MSE loss $\mathcal{L}_{\text{depth}}$ so the global objective is:

$$\mathcal{L} = 0.5\mathcal{L}_{\text{sdf}} + \mathcal{L}_{\text{pbr}} + 0.5\mathcal{L}_{\text{def}} + 0.1\mathcal{L}_{\text{mask}} + 0.1\mathcal{L}_{\text{depth}}.$$

1000 The texture refiner uses only PBR loss and the deferred shading loss: $\mathcal{L}_{\text{pbr}} + 0.5\mathcal{L}_{\text{def}}$.
1001 In each training batch, we randomly sample 4 views per scene as source input views I_i , and another 4
1002 target views I^{tgt} , into which we render the sdf-field predicted by EmuLRM, or the mesh predicted by
1003 the texture refiner. We then evaluate the aforementioned losses in the target views. 3 scenes per GPU
1004 are sampled randomly, and we train on 64 GPUs NVIDIA A100 gpus, yielding an effective batch
1005 size of $3 \times 4 \times 64 = 768$ images. The total loss has been optimized using Adam [38] with learning
1006 rate 10^{-4} for 13K steps.

1007 A.6.3 Deferred shading loss ablation



Figure 11: Generated assets for the prompt: “A cat made of <MATERIAL>”. Emu3D predicts various plausible PBR material maps leading to realistic interaction with the environment light (sphere-mapped in the center)



Figure 12: Using a deferred shading loss on rendered channels enhances PBR quality, resulting in more defined metallicity and roughness, such as increased metallicity in the lantern’s metal parts and decreased roughness in its glass parts.

1008 Besides verifying quantitatively the benefits of the deferred shading loss \mathcal{L}_{def} in Tab. 1, we also
1009 provide a qualitative proof in Fig. 12. Specifically, the PBR materials predicted from albedo&shaded
1010 channels exhibit better metallicity map on the actual metallic parts of the 3D lantern asset.

1011 A.6.4 Direct SDF loss \mathcal{L}_{sdf}

1012 We follow Azinovic *et al.*’s [1] direct SDF supervision for the SDF field. Given a pixel p in an image
1013 and the sampled points S_p on the ray corresponding to the pixel, the direct SDF loss is computed as

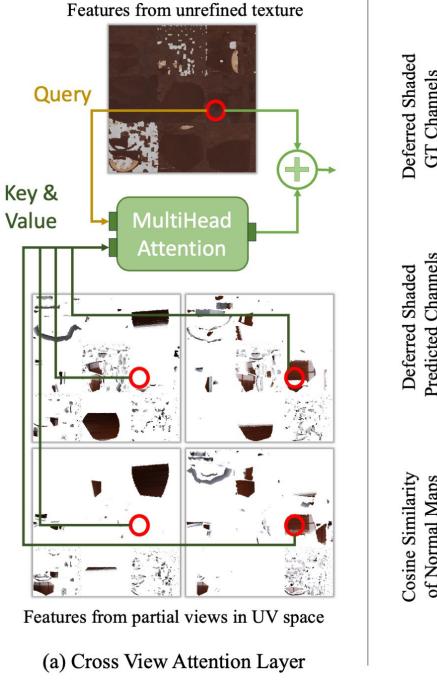
$$\mathcal{L}_{\text{sdf}}(p) = \mathcal{L}_{\text{sdf}}^{\text{tr}}(p) + 0.01\mathcal{L}_{\text{sdf}}^{\text{fs}}(p). \quad (7)$$

1014 $\mathcal{L}_{\text{sdf}}^{\text{fr}}$ is a ‘free-space’ objective, which forces the MLP to predict a value of 1 for samples $s \in S_p^{\text{fs}}$
1015 which lie between the camera origin and the truncation region of a surface:

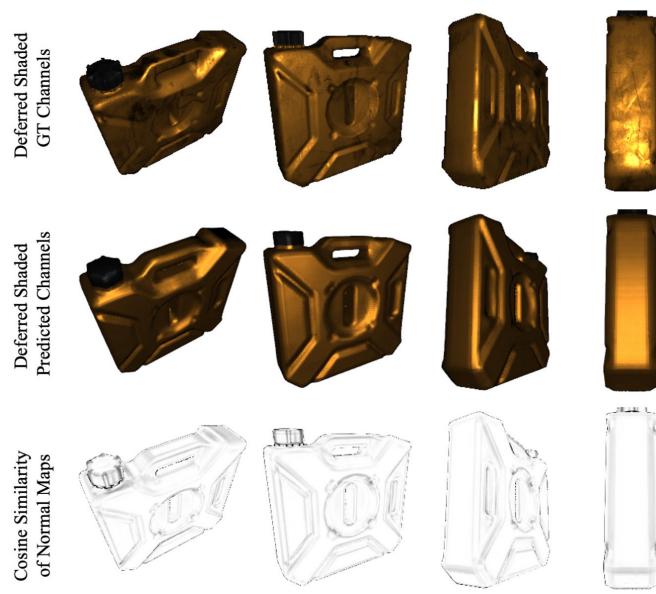
$$\mathcal{L}_{\text{sdf}}^{\text{fr}}(p) = \frac{1}{|S_p^{\text{fr}}|} \sum_{s \in S_p^{\text{fs}}} (D_s - 1)^2 \quad (8)$$

1016 where D_s is the predicted SDF from the MLP. For samples within the truncation region ($s \in S_p^{\text{tr}}$), we
1017 apply $\mathcal{L}_{\text{sdf}}^{\text{tr}}$, the signed distance objective of samples close to the surface.

$$\mathcal{L}_{\text{sdf}}^{\text{tr}}(p) = \frac{1}{|S_p^{\text{tr}}|} \sum_{s \in S_p^{\text{tr}}} (D_s - \hat{D}_s)^2 \quad (9)$$



(a) Cross View Attention Layer



(b) Deferred Shading Loss Example

Figure 13: **(a)** Illustration of Cross-View Attention. Cross-view attention facilitates communication between the UNet branches processing the predicted texture features and the UV space projected input views. This layer blends the predicted texture features with the UV projected input view features based on their match using a multiheaded attention mechanism. **(b)** Example of Deferred Shading Loss Calculation. Deferred shading computes pixel shading using albedo, metallicity, roughness, normals, object position, and light source position. We apply it to both the ground truth channels (top) and the predicted channels (middle). The error is calculated as the difference between the two, weighted by the similarity between ground truth normals and predicted normals, to avoid penalizing shading errors due to incorrect normals.

1018 A naïve PyTorch implementation of this is memory intensive, because of the evaluation of $B \times$
 1019 $H \times W \times N_{\text{ray}}$ points, where B, H, W, N_{ray} are the number of target images in a batch, the height,
 1020 width, and the number of points per ray respectively. Therefore, to support large batch sizes, image
 1021 resolution, and denser point sampling on rays, we implement the direct SDF loss using custom
 1022 Triton [65] kernels.

1023 A.6.5 Depth loss $\mathcal{L}_{\text{depth}}$

The depth loss $\mathcal{L}_{\text{depth}}$ minimizes the mean-squared error between the rendered depth prediction $\mathcal{R}_{\text{depth}}(\cdot | \hat{s}, \pi)$ the ground-truth depth $\mathcal{R}_{\text{depth}}(\cdot | M, \pi)$

$$\mathcal{L}_{\text{depth}} = \left\| \mathcal{R}_{\text{depth}}(\cdot | \hat{s}, \pi) - \mathcal{R}_{\text{depth}}(\cdot | M, \pi) \right\|^2,$$

1024 where $\mathcal{R}_{\text{depth}}(s, \pi)$ is an operator rendering the depth-map of the shape representation s (mesh or an
 1025 SDF) from the viewpoint π .

1026 A.7 Texture refiner

1027 Having described a high-level overview of our texture refiner in Sec. 3.3, here we provide more
 1028 details.

As mentioned, the texture refiner network Φ accepts $N + 1$ texture images K_i in total. The first input to the network is the augmented texture image $K_0 \in \mathbb{R}^{V \times V \times 11}$ given by:

$$\forall v \in [0, V]^2 : K_0(v) = \begin{cases} k(\mathbf{x}_v) \oplus \mathbf{n}(\mathbf{x}_v) \oplus \mathbf{x}_v, & \text{if } v \in \text{Im}(\phi), \\ \mathbf{0}, & \text{otherwise,} \end{cases} \quad \text{where } \mathbf{x}_v = \phi(v).$$

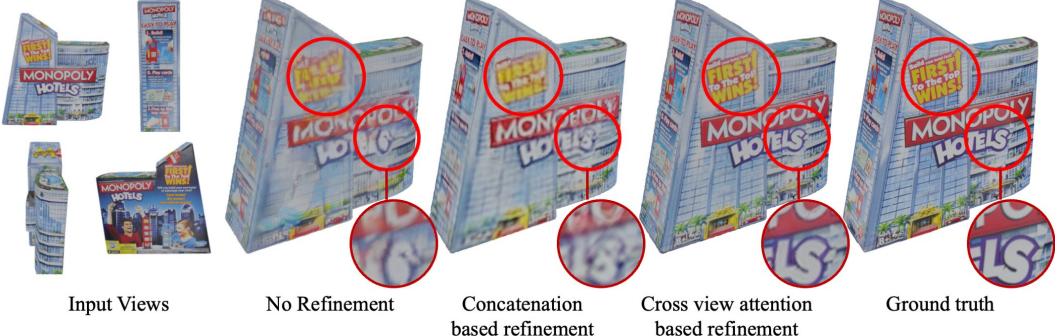


Figure 14: **Texture Refiner Ablation.** When no texture refiner is used, EmuLRM produces textures with washed-out details. Using a refining UNet that concatenates UV space projected input views and predicted texture recovers some details, but our proposed cross-view attention-based refinement is most effective at recovering details and harmonizing the generated texture in unseen views.

1029 The condition $v \in \text{Im}(\phi)$ selects ‘valid’ UV points that correspond to mesh points; and \oplus denotes
 1030 channel-wise concatenation, so that $K_0(v)$ is the stack of the 5 PBR parameters $k(\mathbf{x}_v)$ from EmuLRM,
 1031 normal $\mathbf{n}(\mathbf{x}_v)$, and the 3D point \mathbf{x}_v .

In addition to K_0 , we input to the network Φ texture images K_i , each extracted by looking up information from the corresponding input view I_i directly (thus sidestepping EmuLRM). As noted above, each valid texture point v corresponds to a unique 3D point $\mathbf{x}_v = \phi(v) \in M$ on the mesh, which in turn projects to a pixel $u = \pi_i(\mathbf{x}_v)$ in the image I_i . Let $\chi_i(v) \in \{0, 1\}$ be the flag that tells if point \mathbf{x}_v is *visible* in image I_i or not. When point \mathbf{x}_v is visible in several views, it is best measured in the most frontal one, which is captured by the cosine $\omega_o \cdot \mathbf{n}(\mathbf{x}_v)$ between the normal \mathbf{n} at \mathbf{x}_v and the ray direction $\omega_v \propto \mathbf{x}_0 - \mathbf{x}_v$. All this information is packed into additional texture images $K_i \in \mathbb{R}^{V \times V \times (D+1)}$ by setting:

$$\forall v \in [0, V]^2 : K_i(v) = \begin{cases} I_i(\pi_i(\mathbf{x}_v)) \oplus (\omega_v \cdot \mathbf{n}(\mathbf{x}_v)), & \text{if } v \in \text{Im}(\phi) \text{ and } \chi_i(v) = 1, \\ 0, & \text{otherwise.} \end{cases}$$

1032 The texture network Φ is a U-Net that takes as input the texture augmented texture image K_0 and
 1033 outputs the final enhanced texture $K \in \mathbb{R}^{V \times V \times 5}$. This network also fuses information from the
 1034 view-specific texture images K_i . The goal is to select, for each UV point v , which of the N input
 1035 views provides the best information. This is achieved via cross-attention. Specifically, each K_i
 1036 is processed in parallel by another U-Net, and the first queries information across all the others
 1037 via multi-head cross attention. In Fig. 13 (a), we provide an illustration of the latter cross-view
 1038 attention layer. Additionally, Fig. 14 compares the UNet’s cross-view attention with a simple feature
 1039 concatenation, revealing superiority of the former.

1040 A.8 Physically-Based Rendering: Radiance, BRDFs, and models

1041 We briefly summarise key notion of radiometry and standard BRDF models, and then provide a
 1042 precise expression of the BRDF model used in Emu3D.

1043 A.8.1 Radiance

1044 The *radiant flux* Φ is the electromagnetic power flowing through a particular surface $A \subset \mathbb{R}^3$ oriented
 1045 by the unit normal \mathbf{n} . The *radiance* $L(p, \omega)$ is the radiant flux density at p towards a particular
 1046 direction ω per unit *orthogonal* area dA_{\perp} and per unit solid angle $d\Omega$. The unit vector ω points at
 1047 the direction of propagation of the flux.

1048 The flux density is measured with respect to an area which is orthogonal to the direction of propagation
 1049 ω . In fact, the energy flow is the same through all areas that cut the same ‘tube of flux’; the specific
 1050 area is irrelevant and not a property of the radiation. This dependency is removed by considering the
 1051 normalized area dA_{\perp} , which is orthogonal to the direction of the flux.

Because the radiance is expressed in units of orthogonal area dA_{\perp} , in order to compute the flux through the surface patch dA , which may not be orthogonal, we must account for the *foreshortening factor*, which relates the areas dA and dA_{\perp} :

$$dA_{\perp} = |\langle \mathbf{n}, \omega \rangle| dA.$$

With this, the flux that passes through dA towards direction ω in the solid angle $d\Omega$ is

$$d\Phi = L(p, \omega) |\langle \mathbf{n}, \omega \rangle| dA d\Omega.$$

1052 Note that $d\Phi$ depends on both ω and \mathbf{n} whereas L only depends on ω . This reinforces the notion that
1053 \mathbf{n} is a property of the surface, not of the radiation.

1054 A.9 Reflectance models

1055 Let p be a point on a surface A that separates two media and let dA be a surface patch sitting at p .
1056 Let \mathbf{n} be the normal at this point. We now consider the case where A separates air or empty space
1057 from an opaque object, with \mathbf{n} pointing towards the outside of this object.

1058 Let ω be an orientation on the same side of the surface as \mathbf{n} , i.e., such that $\langle \mathbf{n}, \omega \rangle \geq 0$. From the
1059 viewpoint of the object, we interpret $L(p, \omega)$ as *outgoing radiant flux* and $L(p, -\omega)$ as *incoming
1060 radiant flux*. These two quantities are related by the *Bidirectional Reflectance Distribution Function
1061 (BRDF)* f , defined such that:

$$\frac{dL(p, \omega_o)}{d\Omega_i} = f(p, \omega_i, \omega_o) \langle \mathbf{n}, \omega_i \rangle L(p, -\omega_i). \quad (10)$$

1062 In this definition, for convenience both ω_i and ω_o are taken on the same side as \mathbf{n} (hence the negative
1063 sign in front of ω_i). A useful consequence is that we do not need to take the absolute value of the
1064 inner product $\langle \mathbf{n}, \omega_i \rangle$ as this is positive by definition.

1065 The BRDF thus takes the radiation receives from direction $-\omega_i$ and distributes it along various
1066 outgoing directions ω_o . Integrating over all incoming directions, gives use the overall radiation
1067 reflected towards ω_o :

$$L(p, \omega_o) = \int_{H(\mathbf{n})} \frac{dL(p, \omega_o)}{d\Omega_i} d\omega_i = \int_{H(\mathbf{n})} f(p, \omega_i, \omega_o) \langle \mathbf{n}, \omega_i \rangle L(p, -\omega_i) d\Omega_i. \quad (11)$$

1068 where $H(\mathbf{n}) = \{\omega : \langle \mathbf{n}, \omega \rangle \geq 0\}$ is the hemisphere. Next, we provide common basic models for
1069 the BRDF function in PBR.

1070 A.9.1 Diffuse reflectance

In diffuse reflectance, the radiation is absorbed by the material, internally scattered in random directions, and output again to give rise to a uniform distribution. Namely, the *diffuse BRDF* is:

$$f(p, \omega_i, \omega_o) = \frac{R}{\pi}$$

1071 where $0 \leq R \leq 1$ is the fraction of power reflected by the diffusion process. The $1/\pi$ factor ensures
1072 that the total energy is conserved when $R = 1$.

1073 A.10 Specular reflectance

The reflection for a perfectly flat interface between two media at p is *specular*: the incoming light radiation $-\omega_i$ is partially reflected in the specular direction $\omega_o = r(\mathbf{n}, \omega_i) = 2\mathbf{n}\langle \mathbf{n}, \omega_i \rangle - \omega_i$, and partially transmitted. This phenomena is characterised by *Fresnel's equations*, which are derived from Maxwell's equations, utilising continuity conditions for the electromagnetic field at the interface between the two media. Fresnel's equations describe the planar radiation in full, including its polarisation (in the most general case, using phasors, and thus complex numbers). For graphics, we assume that light is unpolarized, so we only calculate the power. The fraction of power reflected is given by Fresnel's coefficient (using Schlick's approximation [70]):

$$F(\langle \mathbf{n}, \omega_i \rangle) = F_0 + (1 - F_0) |\langle \mathbf{n}, \omega_i \rangle|^5, \quad F_0 = \left(\frac{\hat{n}_1 - \hat{n}_2}{\hat{n}_1 + \hat{n}_2} \right)^2.$$

1074 Here \hat{n}_1 and \hat{n}_2 are the indices of reflectivity of the two media, respectively. This equation is valid
 1075 for dielectrics (non-metallic objects), but also used as an approximation for metals by tweaking F_0 .

In order to write this relation as a BRDF, we write:

$$L(p, \omega_o) = \int_{H(\mathbf{n})} f(p, \omega_i, \omega_o) \langle \mathbf{n}, \omega_i \rangle L(p, -\omega_i) d\Omega_i = R(\langle \mathbf{n}, r(\mathbf{n}, \omega_o) \rangle) L(p, -r(\mathbf{n}, \omega_o)).$$

1076 Hence, the BRDF must be a delta function centered at $\omega_i^* = r(\mathbf{n}, \omega_o)$:

$$f(p, \omega_i, \omega_o) = \frac{F(\langle \mathbf{n}, \omega_o \rangle)}{\langle \mathbf{n}, \omega_o \rangle} \delta_{r(\mathbf{n}, \omega_o)}(\omega_i) \quad (12)$$

1077 where we used the fact that $\langle \mathbf{n}, r(\mathbf{n}, \omega_o) \rangle = \langle \mathbf{n}, \omega_o \rangle$ and where $\delta_{r(\mathbf{n}, \omega_o)}$ is the delta distribution
 1078 centered at $r(\mathbf{n}, \omega_o)$.

1079 A.11 Microfacet models

Rough surfaces can be thought of as a collection of randomly-oriented flat microfacets, each reflecting light in a specular fashion. Consider a point p on a surface and incoming and outgoing radiation directions ω_i and ω_o . If the point contains a microfacet that enables light to reflect from ω_i to ω_o , then the normal of the microfacet must be $\mathbf{m} \propto \omega_i + \omega_o$. If the microfacet is oriented elsewhere, then no light flows in the direction ω_o . Hence, we can write

$$\mathbf{m} = h(\omega_i, \omega_o) = \frac{\omega_i + \omega_o}{|\omega_i + \omega_o|}$$

1080 as a function of the incoming and outgoing radiation. This is called *half vector* as it sits in between
 1081 the two vectors ω_i and ω_o .

Now we wish to derive the macro BRDF $f(p, \omega_i, \omega_o)$ from the micro BRDF $F(p, \omega_i, \omega_o | \mathbf{m})$, where we have emphasized the fact that the BRDF is oriented relative to the microfacet normal \mathbf{m} . A short calculation [87] shows that:

$$f(p, \omega_i, \omega_o) = f_m(p, \omega_i, \omega_o | \mathbf{m}) \frac{\langle \mathbf{m}, \omega_i \rangle}{\langle \mathbf{n}, \omega_i \rangle} \frac{\langle \mathbf{m}, \omega_o \rangle}{\langle \mathbf{n}, \omega_o \rangle} \frac{1}{\langle \mathbf{m}, \mathbf{n} \rangle}.$$

In practice, there is a distribution over possible surface normals \mathbf{m} , characterised by the *microfacet distribution function* $D(\mathbf{m} | \mathbf{n})$. The latter is defined such that $D(\mathbf{m} | \mathbf{n}) dA d\Omega_m$ is the total area of the microfacets within patch dA of the macrosurface with orientation in $d\Omega_m$. In practice, only part of the microfacet is visible and illuminated, depending on the interaction with other facets. This is accounted for by the *shadowing-masking function* $G(\omega_i, \omega_o, \mathbf{m}, \mathbf{n}) \in [0, 1]$. The expected reflectance is thus:

$$f(p, \omega_i, \omega_o) = \int_{H(\mathbf{n})} \frac{\langle \mathbf{m}, \omega_i \rangle}{\langle \mathbf{n}, \omega_i \rangle} \frac{\langle \mathbf{m}, \omega_o \rangle}{\langle \mathbf{n}, \omega_o \rangle} f_m(p, \omega_i, \omega_o | \mathbf{m}) D(\mathbf{m} | \mathbf{n}) G(\omega_i, \omega_o, \mathbf{m}, \mathbf{n}) d\Omega_m.$$

1082 Plugging Eq. (12) in the value for the mirror-like reflectance f_m for each microfacet, we get [87]:

$$f(p, \omega_i, \omega_o) = \frac{F(\langle \mathbf{h}, \omega_o \rangle) D(\mathbf{h} | \mathbf{n}) G(\omega_i, \omega_o, \mathbf{h}, \mathbf{n})}{4 \langle \mathbf{n}, \omega_i \rangle \langle \mathbf{n}, \omega_o \rangle} \quad \text{where } \mathbf{h} = h(\omega_i, \omega_o).$$

1083 A.12 Standard microfacet models

1084 Here, we discuss common choices for the functions D and G in standard PBR models.

The Torrance-Sparrow model. The oldest such model is due to Torrance and Sparrow [84]. They simply assume a Gaussian model for the microfacet distribution function:

$$D(\mathbf{m} | \mathbf{n}) = b \exp(-\alpha^2 \theta), \quad \text{where } \cos \theta = \langle \mathbf{m}, \mathbf{n} \rangle,$$

and b is a suitable normalization constant. For the shadowing-masking function they pick:

$$G(\omega_i, \omega_o, \mathbf{m}, \mathbf{n}) = \min \left\{ 1, \frac{2 \langle \mathbf{m}, \mathbf{n} \rangle \langle \mathbf{n}, \omega_i \rangle}{\langle \mathbf{m}, \omega_i \rangle}, \frac{2 \langle \mathbf{m}, \mathbf{n} \rangle \langle \mathbf{n}, \omega_o \rangle}{\langle \mathbf{m}, \omega_o \rangle} \right\}.$$

1085 This function has a simple geometric derivation under a basic geometric model of the microfacets.

The Beckmann-Spizzichino-Smith model. Beckmann and Spizzichino [2] suggested the model:

$$D(\mathbf{m}|\mathbf{n}) = \frac{1}{\pi\alpha^2 \cos^4 \theta} e^{-\frac{\tan^2 \theta}{\alpha^2}} \quad \text{where } \cos \theta = \langle \mathbf{m}, \mathbf{n} \rangle.$$

Smith [74] noted that the shadowing-masking function should be derived from the microfacet distribution function, which describes the micro-geometry of the surface. They also proposed a factorized model $G(\omega_i, \omega_o, \mathbf{m}, \mathbf{n}) = G_1(\omega_i, \mathbf{n})G_1(\omega_o, \mathbf{n})$. For Beckmann's distribution, the Smith shadowing-masking function is given by:

$$G_1(\omega, \mathbf{n}) = \frac{2}{1 + \operatorname{erf}(a) + \frac{1}{a\sqrt{\pi}} e^{-a^2}} \quad \text{where } a = \frac{1}{\alpha \tan \theta_\omega} \text{ and } \cos \theta_\omega = \langle \mathbf{n}, \omega \rangle.$$

1086 **The GGX model.** The GGX model by [87] is a variant of the Beckmann model, with slightly
1087 different microfacet distribution and shadowing-masking function:

$$D(\mathbf{m}|\mathbf{n}) = \frac{\alpha^2}{\pi \cos^4 \theta (\alpha^2 + \tan^2 \theta)^2}, \quad G_1(\omega, \mathbf{n}) = \frac{2}{1 + \sqrt{1 + \alpha^2 \tan^2 \theta_\omega}}, \quad (13)$$

1088 where $\cos \theta = \langle \mathbf{m}, \mathbf{n} \rangle$ and $\cos \theta_\omega = \langle \mathbf{n}, \omega \rangle$.

1089 A.12.1 The BRDF model used in Emu3D

The BRDF model used in our paper combines diffuse and GGX BRDFs:

$$f(\omega_i, \omega_o | k(\mathbf{x}), \mathbf{n}) = \frac{R}{\pi} + \frac{F(\mathbf{h}|\mathbf{n})D(\mathbf{h}|\mathbf{n})G_1(\mathbf{n}, \omega_i)G_1(\mathbf{n}, \omega_o)}{4(\mathbf{n} \cdot \omega_i)(\mathbf{n} \cdot \omega_o)}.$$

The first term is the *diffuse component* (Lambertian reflection), where $R \in [0, 1]$ is the fraction of light power reflected by diffusion. The second term in the *specular component*, where F is Schlick's approximation [70] $F(\omega_i|\mathbf{n}) = F_0 + (1 - F_0)(1 - \mathbf{n} \cdot \omega_i)^5$ of Fresnel's reflectance and where $F_0 \in [0, 1]$ is the Fresnel coefficient at normal incidence. The unit vector $\mathbf{h} \propto \omega_i + \omega_o$ is the *half vector*, which is the orientation needed to reflect ω_i into ω_o by the rough material's microfacets (generally different from but averaging to \mathbf{n}). The function D and G_1 are the microfacet distribution function and the shadowing-masking function given by:

$$D(\mathbf{m}|\mathbf{n}) = \frac{\alpha^2}{\pi((\mathbf{m} \cdot \mathbf{n})^2(\alpha^2 - 1) + 1)^2}, \quad G_1(\mathbf{n}, \omega) = \frac{2(\mathbf{n} \cdot \omega)}{\mathbf{n} \cdot \omega + \sqrt{(\mathbf{n} \cdot \omega)^2(\alpha^2 - 1) + \alpha^2}}.$$

These are the same as Eq. (13) with the trigonometric functions expanded in terms of dot products. In this model, the reflectance R and reflectance at normal incidence F_0 are RGB triplets. The specular highlight color F_0 is approximately white (equal to 1) for dielectrics, and colored for metals; furthermore, metals have no diffuse component ($R = 0$). Thus, we introduce the parameter *metallicity* $\gamma \in [0, 1]$ and the *base color* $\rho_0 \in [0, 1]^3$ and define:

$$R = \rho_0(1 - \gamma), \quad F_0 = \mathbf{1}(1 - \gamma) + \rho_0\gamma.$$

1090 In this manner, the albedo is used as diffuse color or specular color depending on whether the material
1091 is a dielectric or a metal. In the paper, we call the parameter ρ_0 *albedo* as it is a better known concept;
1092 however, in this model ρ_0 is the albedo only when $\gamma = 0$, i.e., when the object is a perfect dielectric.

1093 The BRDF is thus fully described by the albedo ρ_0 , the roughness α , and the metallicity γ , for a total
1094 of five scalar parameters. Hence, the LRM predicts the triplet $k(\mathbf{x}) = (\rho_0, \gamma, \alpha)$ at each 3D point \mathbf{x} .