
Diffusion-based photo-realistic rendering with scalable synthetic-real paired data

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Photorealistic rendering aims to produce images indistinguishable from real-world
2 photographs. Traditional rendering techniques, while effective, rely heavily on syn-
3 thetic models with intricate material properties. While neural rendering methods
4 offer a potential solution to this problem, they often necessitate data from costly
5 capturing equipment like the Light Stage or resort to low-quality synthetic data,
6 hindering their ability to achieve photo-realistic rendering. To address these chal-
7 lenges, we propose a novel image rendering framework including a data generation
8 method and a neural rendering model. Our data generation method can create
9 synthetic-real data pairs using intrinsic decomposition methods, leveraging intrinsic
10 images similar to G-buffers in the traditional rendering pipeline. Additionally,
11 we introduce a photorealistic image synthesis method based on diffusion models,
12 enhancing the generalization capabilities of our framework. This framework allows
13 for scalable data generation and photorealistic rendering for low-quality synthetic
14 objects. Experiments show that our method can not only render comparable images
15 with sophisticated synthetic 3D models but can fulfill state-of-the-art rendering for
16 low-quality synthetic 3D models.

17 1 Introduction

18 In the pursuit of photorealistic rendering, the goal is to generate images that are virtually indistinguish-
19 able from real-world photographs, a necessity in fields like game production and immersive virtual
20 reality. With the evolution of computer hardware, the advent of physically based rendering (PBR)
21 utilizing recursive path tracing has become feasible. Over recent decades, researchers have dedicated
22 significant efforts to crafting sophisticated rendering models[41][53][5] that meticulously account
23 for lighting, materials, and object geometry to achieve photorealism. These rendering techniques
24 are now commonplace in production engines such as Unreal Engine[15] and Unity[52]. However,
25 their efficacy often hinges on the availability of high-quality CAD models with intricate material
26 properties, necessitating extensive manual labor. Consequently, rendering photorealistic images
27 remains a challenge when dealing with low-quality CAD models generated through manual design or
28 techniques like Multi-View-Stereo (MVS).

29 With the advent of deep learning, the prospect of generating photorealistic images using end-to-
30 end neural networks has become a reality. Previous research[51][55][40] has utilized reflectance
31 fields obtained from Light Stage[11][12] captures to construct datasets, which are then employed
32 to train rendering networks. Leveraging these high-quality datasets, neural networks have achieved
33 state-of-the-art results. However, acquiring such datasets poses a challenge for consumers, as it
34 typically requires access to expensive Light Stage equipment. To circumvent this limitation, some
35 methodologies[32][70][61] have emerged to synthesize training datasets. While these approaches
36 prove effective to a certain extent, networks trained solely on synthetic datasets often struggle to

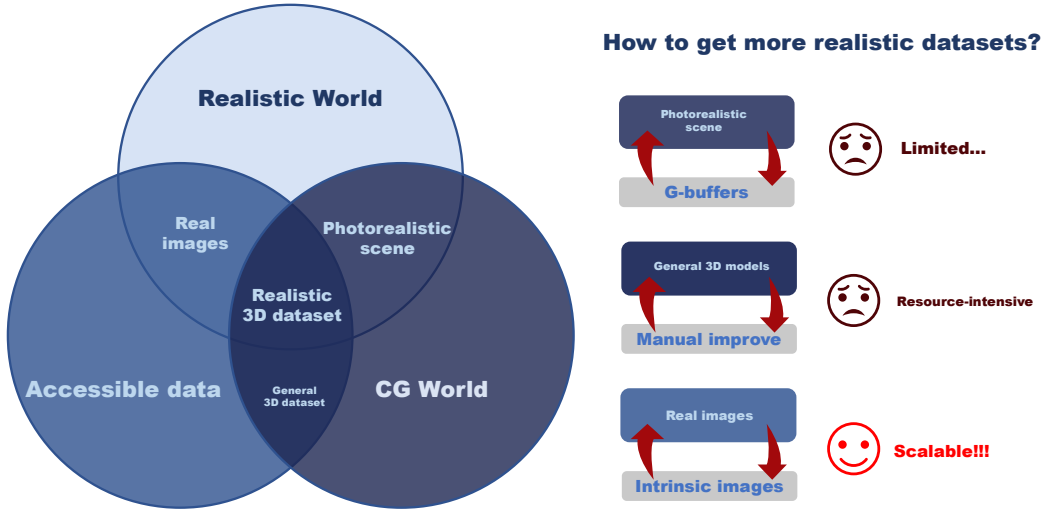


Figure 1: The teaser shows our major contribution to solving the problem of lacking realistic paired images for neural rendering. According to the figure on the left side, we can enrich the realistic paired images from three different resources - Photorealistic 3D scenes, general 3D datasets, and real images. Though there is a large amount of photorealistic 3D scenes in film productions that can be used to render high-quality paired datasets, it is difficult to get access to them. While manually creating more realistic data may seem intuitive, it is a resource-intensive process, demanding significant time and financial investment. Numerous general 3D datasets exist, but their quality may not be as satisfactory. Compared to these two resources, generating realistic paired data from real images is promising since the novel framework we proposed, based on intrinsic images and the neural rendering method, has proved to be practical.

37 generalize well to real-world scenarios, hindering their ability to render photorealistic images for
 38 synthetic objects.

39 In response to the challenge of limited paired data for rendering synthetic objects with photorealism,
 40 we present a novel data generation pipeline capable of producing synthetic-real data pairs. Our
 41 key insight stems from the observation that intrinsic images derived from intrinsic decomposition
 42 methods exhibit a domain similarity to the G-buffers within the computer graphics (CG) rendering
 43 pipeline, as shown in Figure 2. Specifically, we extract irradiance images, specular shading images,
 44 and albedo images as the source data, as they encompass low-frequency lighting information, high-
 45 frequency view-dependent information, and identity information, respectively. To obtain these images
 46 from natural scenes, we employ the IntrinsicAnything[9] method which is trained on the synthetic
 47 Objaverse dataset[13]. While IntrinsicAnything excels in generating high-quality albedo images
 48 and specular shading images, it lacks the capability to produce irradiance images. To address this
 49 limitation, we introduce a novel method designed to generate all aforementioned intrinsic components.

50 With the generated data, we further propose a photorealistic image synthesis method based on the
 51 diffusion model. Recently, the diffusion model[19][46] has demonstrated state-of-the-art performance
 52 across various tasks[24][36][33][16]. These methodologies have shown that diffusion models, even
 53 when fine-tuned on small datasets, exhibit strong generalization capabilities. This is primarily
 54 attributed to the robust priors learned by diffusion models from extensive real data. In this paper,
 55 we also capitalize on the strength of the diffusion prior. Specifically, we first use pre-trained visual
 56 encoders to extract global and local features of the intrinsic images. Then these features are injected
 57 into a pre-trained text-to-image diffusion model as guidance.

58 Through our data generation pipeline, we can readily scale up the training dataset by gathering
 59 extensive natural images from the internet and processing them. This scalability is pivotal for
 60 enhancing the generalization capabilities of neural networks. Moreover, leveraging the robust
 61 diffusion prior, our fine-tuned diffusion model excels not only in rendering comparable images for
 62 high-quality CAD models but also in generating photorealistic images for even the most low-quality
 63 synthetic objects. In summary, our main contributions are:



Figure 2: Visualization of the G-buffers and intrinsic images.

- 64 • We propose a novel training framework for photo-realistic image synthesis which includes a
- 65 data generation method and a neural rendering model.
- 66 • We achieve state-of-the-art photorealistic rendering results.

67 2 Related work

68 2.1 Photorealistic Rendering

69 Modern graphics heavily rely on Physically Based Rendering (PBR) techniques to achieve pho-
70 torealism. PBR methods[41][4][53][5], incorporating lighting, materials, and geometry, render
71 meticulously crafted CAD models according to the rendering equation. Early graphic production
72 saw the prominence of models like the Phong[41] and its variant, the Blinn-Phong[4] model, which
73 delivered a commendable performance. However, sophisticated material properties were necessitated
74 for more realistic rendering, leading to the development of advanced models like the Cook-Torance
75 GGX[53] and Disney GGX[5]. While these PBR methods enable photo-realistic rendering, they
76 demand high-quality CAD models with intricate materials and geometry, incurring substantial costs.
77 To address this challenge, Paul Debevec *et al.* introduced the Light Stage[12], enabling the capture
78 of reflectance fields for human portraits, facilitating rendering under any natural environment[10].
79 Although combining reflectance fields with environment maps can produce nearly natural images,
80 this approach is limited to subjects with pre-captured reflectance fields. Moreover, accessibility to
81 such technology remains restricted, constraining its widespread application.

82 Another significant avenue of photorealistic rendering is neural rendering, which harnesses the
83 power of neural networks to generate high-quality images. A common task in neural rendering is
84 relighting, wherein the goal is to alter the appearance of subjects in an image to match a different
85 target environment light. Existing learning-based approaches either[51][55][40][43][67][25] utilize
86 images generated from reflectance fields[11] or synthesized datasets[32][70][61] to train end-to-end
87 networks for single image relighting. Some methods[51][70][43][67] alter the light information and
88 re-render input images fully in the latent space, while others[32][55][40][61] first estimate intrinsic
89 image properties before combining them with target lighting for neural rendering. SIPR[51] and
90 DSIPR[70] inject target lighting into the network bottleneck and use the decoder for rendering.
91 PhotoApp[43] edits lighting in the latent space of StyleGAN[23], then decodes the latent code for
92 rendering. NVPR[67] introduces self-supervised losses to disentangle lighting and identity features,
93 which helps to render with novel light. Li et al.[32] design a multi-bounce scheme to ease the problem
94 of not considering the global illumination and a cascade framework for narrowing the errors of the
95 predictions. Thanks to the existence of the shadow maps and specular maps in the generated dataset,
96 Wang et al. [55] design the SS network to predict the specular maps and shadows maps which are
97 input into the composition network together with light and albedo to composite a realistic rendering
98 image. Total Relighting[40] utilizes convolved environment maps and normal maps to generate
99 light maps containing explicit lighting cues. These light maps, along with other intrinsic properties,
100 serve as inputs to a shading network for rendering photorealistic images. Due to the absence of
101 publicly available light stage datasets, Lumos[61] relies on purchased 3D face scans to generate large-
102 scale relighting pairs. Following the Total Relighting framework, Lumos initially trains networks
103 using synthetic datasets and then refines the results by learning a residual map that minimizes the
104 domain gap between synthetic and real albedo using a real dataset. By leveraging extracted image

105 intrinsics, SwitchLight[25] combines the Cook-Torrance model for initial relighting with a neural
106 network for enhanced refinement. Recently, diffusion models[19][46] trained on large-scale natural
107 images have demonstrated impressive performance across various vision tasks. To fully exploit the
108 priors embedded in diffusion models for rendering, several recent works[14][42][27][63] have made
109 notable attempts. DiffusionRig[14] and DiFaReli[42] first utilize DECA to predict FLAME[30]
110 parameters and spherical harmonics (SH) light information. Subsequently, they render physical
111 buffers containing light information to guide the diffusion process. LightIt[27] predicts shading maps
112 for outdoor scenes, utilizing them as conditions for generating relighted images. DiLightNet[63], by
113 pre-defining multiple roughness levels, generates radiance hints as guidance for the diffusion model.
114 In our paper, we aim to condition the diffusion model on albedo, irradiance, and specular shading
115 maps to achieve photo-realistic rendering.

116 2.2 Inverse Rendering

117 Inverse rendering endeavors to recover the intrinsic properties (such as geometry, materi-
118 als, and lighting) of a scene or object, either in 2D[17][28][20][31][62][59][35][58][56] or
119 3D[7][57][1][54][66][50][68][65][69][29][21][64][34][38] space. In 3D inverse rendering, some
120 methods[7][1][54][57] address the problem based on explicit mesh. DIB-R++[7] first predicts the
121 mesh with UV mapping from the input single image, then optimizes the material textures and
122 lighting via the differentiable renderer. Azinovic *et al.*[1] and SunStage[54] estimate the material
123 textures based on the reconstructed FLAME mesh. FIPT[57] calculates the pre-baked shading
124 maps with the provided or optimized mesh and uses the shading maps for materials estimation.
125 Others[50][68][29][66][65][69][34][38] optimize all the intrinsics in implicit neural radiance field.
126 To adapt Nerf[37] for inverse rendering, NERF[50], NerFactor[68] and NEROIC[29] assign materials
127 for each sample point and calculate the outgoing radiance of them with physically-based rendering
128 (PBR). These models[66][65][69][34][38] represent geometry as signed distance function (SDF) and
129 optimize the neural materials of the implicit surfaces.

130 While inverse rendering in 3D has achieved promising results, it often suffers from computational
131 inefficiencies due to the utilization of computation-intensive Multi-Layer Perceptron (MLP). To
132 mitigate this issue, single-image inverse rendering focuses on recovering intrinsic properties using
133 more generalizable neural networks. David *et al.* [17] treats the albedo estimation task as a light
134 diffusion process and iteratively diffuses the image to get albedo. IID[28], trained on indoor synthetic
135 datasets, utilizes diffusion priors to estimate material maps. Due to the lack of ground truth material
136 maps, unsupervised methods[20][31][62][59][35][58][56] estimate the intrinsics based on hand-
137 crafted priors.

138 3 Method

139 Given G-buffers or intrinsic images of an object, our task is to render photo-realistic images. Figure 1
140 shows the overview of our framework. Our innovation is in a novel rendering framework that includes
141 a novel synthetic-real paired data generation method (Section 3.1) and a rendering method leveraging
142 strong diffusion priors. To render photo-realistic images with G-buffers or intrinsic images, our core
143 idea is to extract the global identity features (Section 3.2.1) and the local detailed features (Section
144 3.2.2), then inject (Section 3.3) them into a pre-trained text-to-image diffusion model.

145 3.1 Data generation

146 Though neural rendering shows great performance in fulfilling photo-realism, the data is hard to
147 collect. To solve this problem, we are the first to generate intrinsic-real data pairs for training. Our
148 key observation is that the intrinsic images estimated by intrinsic decomposition methods share a
149 similar domain to the G-buffers of synthetic CAD models. Thus, we can treat the intrinsic domain as
150 a proxy domain, the images from which can be rendered or mapped to natural images. After training,
151 we can obtain photo-realistic images from the G-buffers of the synthetic CAD models.

152 To get the intrinsic images of natural images, we take the state-of-the-art intrinsic decomposition
153 method IntrinsicAnything[9] as the baseline, which can generate the albedo image and specular
154 shading image from a natural image. However, these two images cannot fully represent the informa-
155 tion present in an image, as they only contain the identity and high-frequency lighting information

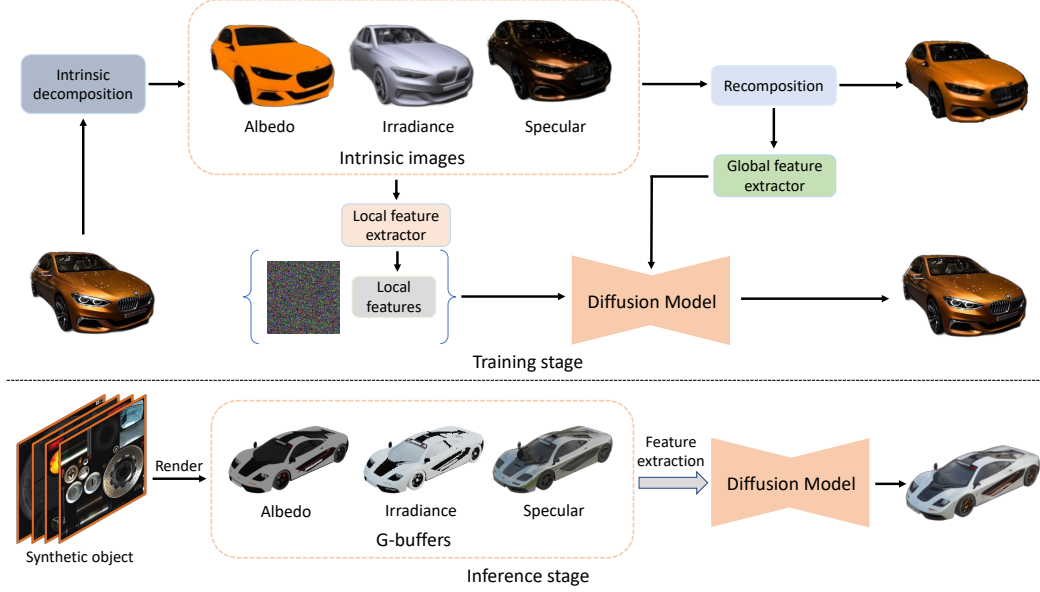


Figure 3: **Overview of our framework.** Given natural images, we first generate intrinsic-real data pairs for training. At the training stage, we extract the global features by encoding the recomposed image into a pre-trained self-supervised visual encoder and extract the local features with a pre-trained VAE encoder. We then inject the features into a pre-trained text-to-image diffusion model for rendering. At the inference stage, we extract features of the rendered G-buffers of the synthetic CAD model, and then inject them into the model for photo-realistic rendering.

156 respectively. To represent the low-frequency lighting information, we need another component, i.e.,
 157 the irradiance image. So the forward pass of intrinsic image generation is defined as:

$$A, S, I = F_{\theta}(I) \quad (1)$$

158 where A , S and I are the albedo image, specular shading image, and irradiance image respectively,
 159 F is the intrinsic decomposition model and θ is the model parameter.

160 To adapt IntrinsicAnything for irradiance estimation, we first render irradiance images for Objaverse
 161 dataset, then we train another diffusion model using these images for irradiance estimation.

162 3.2 Feature extraction

163 Previous visual encoders such as CLIP[44] and Dino[6] are built on natural images. Directly applying
 164 them on intrinsic images may fail to extract representative global features. To address this problem,
 165 we propose to recompose the intrinsic images into a natural image according to the rendering equation.
 166 Then we use the self-supervised model DINOv2[39] to extract the global identity features. As for
 167 the local detailed features, we use the pre-trained VAE encoder to extract the pixel-aligned detail
 168 features.

169 3.2.1 Global feature extraction

170 **Image recomposition.** The rendering equation can be represented as:

$$L = L_d + L_s \quad (2)$$

171 where L_d is the diffuse shading and L_s is the specular shading. The diffuse shading can further be
 172 represented as:

$$L_d = A * I \quad (3)$$

173 where A is the albedo image and I is the irradiance image. So we can get the recomposed image as:

$$R(x) = A(x) * I(x) + S(x) \quad (4)$$

174 where x is the pixel coordinate of the image, $A(x)$, $I(x)$ and $S(x)$ are the albedo image, irradiance
 175 image, and specular shading image at x respectively. Firstly, we transform these images into linear
 176 space. Then we compose them according to Equation 4. Finally, we transform the linear-space image
 177 back to sRGB space by gamma correction. Though the recomposed image loses some high-dynamic-
 178 range information compared to the source natural image due to truncation, we find that the extracted
 179 features are robust enough for rendering.

180 **Feature extraction.** Previous works[49][60] use the CLIP image encoder to extract the global
 181 features. However, the CLIP image embedding is aligned with the text embedding. Text, being a
 182 coarse description, is often insufficient to represent intricate details. Inspired by Anydoor[8], we
 183 choose the Dinov2 as our global feature extractor which encodes the image as a global token and a
 184 patch token. Following AnyDoor, we concatenate the two tokens and use a linear layer to project the
 185 tokens to the diffusion embedding space. The process is defined as:

$$t_l, t_g = F_G(R) \quad (5)$$

$$f_g = L(Con(t_l, t_g)) \quad (6)$$

186 where R is the recomposed image, F_G is the Dinov2 model, t_l is the patch token, t_g is the global
 187 token, Con means concatenation along channel dimension, L is the linear projection layer and f_g is
 188 the global feature in the diffusion embedding space.

189 3.2.2 Local feature extraction

190 Global features encode the identity of the object, while local features encode the details of the
 191 object. Though the VAE model of Stable Diffusion is also trained on natural images, these
 192 methods[24][16][36] proved that they can extract features rich in detailed information for intrinsic
 193 images such as the normal map and the depth map. Thus, we use the pre-trained VAE encoder to
 194 extract the pixel-aligned detail features. Specifically, we first use the VAE encoder to encode the
 195 three intrinsic images. Then we concatenate the three feature maps along the channel dimension. The
 196 concatenated feature map will further be concatenated with the noise image as the input to the Unet
 197 encoder.

198 The process is defined as:

$$f_a, f_i, f_s = F_L(A, I, S) \quad (7)$$

$$f_l = Con(f_a, f_i, f_s) \quad (8)$$

199 where A , I , and S are the albedo image, irradiance image, and specular shading image respectively,
 200 F_L is the pre-trained VAE encoder, f_a , f_i and f_s are the albedo feature map, irradiance feature
 201 map and specular shading feature map respectively, and Con means concatenation along channel
 202 dimension.

203 3.3 Feature injection

204 After obtaining the global and local features, we inject them into a pre-trained text-to-image diffusion
 205 model. In our paper, we take the Stable Diffusion[47] as the backbone which has been demonstrated
 206 robust. We denote the diffusion model as ϵ_θ , so the training objective can be defined as:

$$\mathbb{E}_{t, \mathbf{x}_0, \epsilon, c} [\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t, c)\|_2] \quad (9)$$

207 where $\bar{\alpha}_t$ is the schedule parameter, \mathbf{x}_0 is the target image, ϵ is the noise image sampled from $N(0, I)$,
 208 t is the time step, and c is the extracted features. The training objective is to minimize the L2 distance
 209 between the predicted noise vector and the ground-truth sampled noise vector.

210 Specifically, we inject the global features into the diffusion model by replacing the text embedding
 211 with them and performing cross-attention between them and Unet intermediate features. As for the
 212 local detailed features, we concatenate them with the noise image as the input to the Unet encoder.

213 4 Experiments

214 4.1 Implement Details

215 We implement all our methods in Pytorch. We use the Adam[26] optimizer with an initial learning
216 rate of $1e^{-5}$ and linearly decay the learning rate to 0. During each iteration, we take a batch size of
217 128 for training. We optimize our model for 20k iterations which takes 4 days on 8 Nvidia A800
218 (80GB) GPUs. Other hyperparameters for our network follow the default settings in the Stable
219 Diffusion model.

220 4.2 Dataset and Metrics

Table 1: Comparison of FID, KID, and Inception Score on face and car datasets

Method	Face Dataset			Car Dataset		
	FID ↓	KID×10 ³ ↓	Inception Score ↑	FID ↓	KID×10 ³ ↓	Inception Score ↑
PBR	42.5	37.1±1.2	3.88±0.19	32.1	9.64±0.09	2.25±0.19
Ours	31.1	25.0±1.2	4.18±0.18	28.9	7.17±0.06	2.23±0.17

221 **Datasets** As mentioned, our training framework can generate synthetic-real paired data from real
222 data. To evaluate the performance of the framework and the diffusion-based synthesis method, we
223 conducted experiments on two real datasets: the real car dataset and the real face dataset. We collected
224 the car dataset from the Internet. This dataset contains about 1000 cars, each of which has about
225 6 multi-view images. For the face dataset, we used the FFHQ dataset [22] that provides 70,000
226 high-quality face images in total. We split the car and face datasets so that ninety percent of them are
227 used for training and the rest are used for evaluation. To prove the capability of generalization, we
228 test our model on limited-quality synthetic face and car models collected from TURBOSQUID and
229 CGTrader website. The links of these models are presented in the supplement material. Furthermore,
230 we also evaluate our model on a high-quality synthetic dataset—hyperSim [45] that contains over
231 77400 photorealistic images of 461 indoor scenes with sophisticated and accurate lighting, materials,
232 and geometry.

233 **Metrics** Metrics for evaluating the realism of the generated images have been widely used in the
234 field of generative models. The most common metrics are Inception Score (IS) [48], Fréchet Inception
235 Distance (FID) [18], and Kernel Inception Distance (KID) [2]. Inception score aims to measure the
236 diversity and recognizability of generated images. It uses a pre-trained Inception network to classify
237 generated images and calculates the score based on the entropy of the predicted class distribution. A
238 high IS indicates that the generated images are both diverse and confidently classified into specific
239 categories. However, IS does not consider the distribution difference between the generated images
240 and the real images. FID assesses the similarity between the distributions of generated images and real
241 images. It calculates the Fréchet distance between the feature vectors of real and generated images,
242 extracted from a specific layer of a pre-trained Inception network. Lower FID values indicate that
243 the generated images are closer in distribution to the real images. KID also evaluates the similarity
244 between generated and real images but emphasizes semantic content. It uses a polynomial kernel to
245 compute the Maximum Mean Discrepancy (MMD) between the feature representations of generated
246 and real images, extracted from the third pooling layer (pool3) of a pre-trained Inception network.
247 Lower KID values indicate better alignment between the distributions of generated and real images. It
248 is widely accepted that these three metrics are complementary and should be used together to provide
249 a comprehensive evaluation of the generated images.

250 Table 1 compares FID, KID, and Inception Score across different methods on the FFHQ dataset
251 and our collected car dataset. Please note that PBR means directly computing the rendered results
252 according to the recomposition method explained in Section 3.2.1 and KID has been shown with
253 ×1000 for better readability. The gray part indicates the standard deviation of the metrics.



Figure 4: This figure shows some results of our diffusion-based method compared to the physically based rendering method on synthetic data. The models produced the face-related results trained on the FFHQ [22] dataset and the car-related results trained on our collected car dataset. Please note that face data have slight metallicity due to the capture process, which will result in the irregular specular effect present in the PBR-rendered images. However, our method can still generate realistic results.

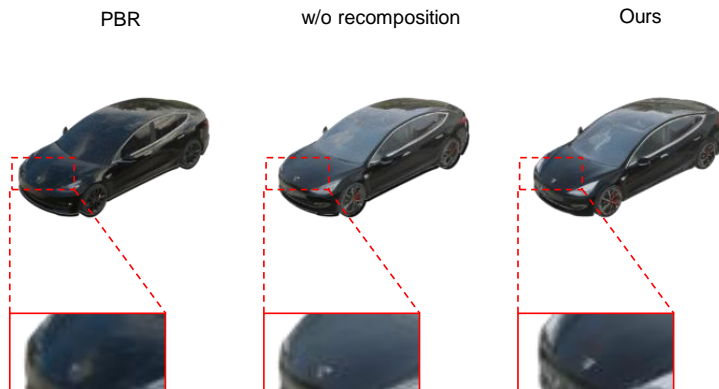


Figure 5: Ablation study for recomposition method.

254 4.3 Comparison with Baselines

255 Since our task aims to render photo-realistic images from G-buffers or intrinsic images, we compare
 256 our method with the physically based rendering (PBR) method. Specifically, we use the Principled
 257 BSDF node in Blender[3] as the description of the object’s material and use environment maps as the
 258 light conditions for rendering. Our results show that we can achieve indistinguishable photo-realistic
 259 results compared to the baseline and even limited automatic realistic relighting. Some of our results
 260 are shown in Figure 4. Other test results on hyperSim[45] dataset are shown in the supplementary
 261 material, which are from .

262 4.4 Ablation Study

263 **Recomposition** The diffusion-based synthesis method is the core of our framework, and we
 264 conduct ablation studies to verify the effectiveness of the recomposition method. It is proved that
 265 the recomposition method can better guide the diffusion model. By recomposition, the generated
 266 images are more realistic and have better details. The examples are shown in the Figure 5. Due to the
 267 low-resolution of our output, the results may not be significantly different. However, the enlarged
 268 areas show kinds of details, which proves the effectiveness of our recomposition method.

269 **Neural rendering method** Before the hyper enthusiasm of diffusion-based image synthesis, there
 270 were also some neural rendering methods based on convolution networks. We also conduct this
 271 method on the above two datasets. The network used here is the Resnet Generator from the
 272 CycleGAN[71] method. The results in Figure 6 show that the results of face data are well sat-
 273 isfactory but those of cars are short of details and realism. Predictably, the network could perform
 274 well with enough data (the FFHQ dataset has nearly 70000 images, and the real car dataset only has
 275 6000 images), which also means lots of data-driven methods could be progressed with our scalable
 276 data generation method.

277 5 Limitations

278 Our method relies highly on the quality of the generated data which is not guaranteed. In practice, we
 279 can find noisy data generated by the intrinsic decomposition method. Using these data for training
 280 may degrade the performance of the model. Meanwhile, the resolution of the generated intrinsic
 281 images is low due to the limitation of the intrinsic model, preventing the attainment of high-quality

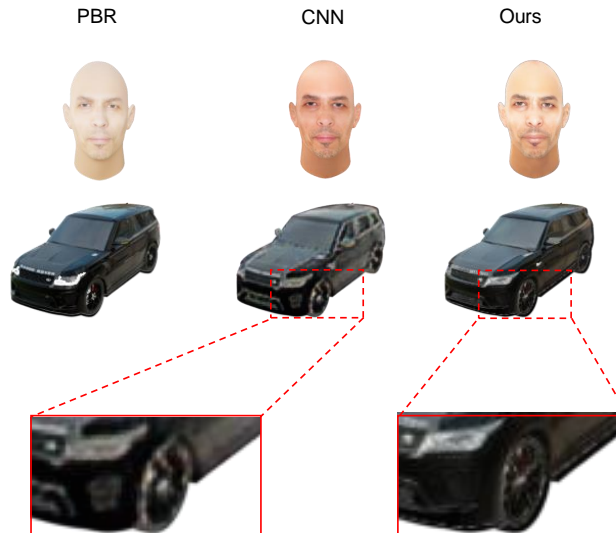


Figure 6: Ablation study for difference neural rendering methods.

282 results. This remains a future work for the intrinsic decomposition task. Besides, our method is based
 283 on the diffusion model, which is computationally expensive and not suitable for real-time rendering.

284 6 Conclusion

285 In this paper, we introduce a novel framework for photo-realistic rendering. This framework includes
 286 a novel synthetic-real paired data generation method and a diffusion-based neural rendering method.
 287 The data generation method leverages the intrinsic decomposition method to generate intrinsic images
 288 for real data. With these data, we extract their global and local features, and then inject them into a
 289 pre-trained text-to-image diffusion model. Qualitative and quantitative experiments demonstrate the
 290 effectiveness of our framework in generating photo-realistic images.

291 References

- 292 [1] Dejan Azinović, Olivier Maury, Christophe Hery, Matthias Nießner, and Justus Thies. High-res
 293 facial appearance capture from polarized smartphone images. In *Proceedings of the IEEE/CVF*
 294 *Conference on Computer Vision and Pattern Recognition*, pages 16836–16846, 2023.
- 295 [2] Marcin Binkowski, Dougal J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying
 296 mmd gans. *arXiv preprint arXiv:1801.01401*, 2018.
- 297 [3] Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender
 298 Foundation, Blender Institute, Amsterdam, 2023.
- 299 [4] James F. Blinn. Models of light reflection for computer synthesized pictures. *SIGGRAPH*
 300 *Comput. Graph.*, 11(2):192–198, jul 1977.
- 301 [5] Brent Burley and Walt Disney Animation Studios. Physically-based shading at disney. In *Acm*
 302 *Siggraph*, volume 2012, pages 1–7. vol. 2012, 2012.
- 303 [6] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski,
 304 and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings*
 305 *of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021.
- 306 [7] Wenzheng Chen, Joey Litalien, Jun Gao, Zian Wang, Clement Fuji Tsang, Sameh Khamis,
 307 Or Litany, and Sanja Fidler. Dib-r++: learning to predict lighting and material with a hybrid

- 308 differentiable renderer. *Advances in Neural Information Processing Systems*, 34:22834–22848,
309 2021.
- 310 [8] Xi Chen, Lianghua Huang, Yu Liu, Yujun Shen, Deli Zhao, and Hengshuang Zhao. Anydoor:
311 Zero-shot object-level image customization. *arXiv preprint arXiv:2307.09481*, 2023.
- 312 [9] Xi Chen, Sida Peng, Dongchen Yang, Yuan Liu, Bowen Pan, Chengfei Lv, and Xiaowei Zhou.
313 Intrinsicanything: Learning diffusion priors for inverse rendering under unknown illumination,
314 2024.
- 315 [10] Paul Debevec. Rendering synthetic objects into real scenes: bridging traditional and image-
316 based graphics with global illumination and high dynamic range photography. In *Proceedings*
317 *of the 25th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH
318 '98, page 189–198, New York, NY, USA, 1998. Association for Computing Machinery.
- 319 [11] Paul Debevec. The light stages and their applications to photoreal digital actors. *SIGGRAPH*
320 *Asia*, 2(4):1–6, 2012.
- 321 [12] Paul Debevec, Tim Hawkins, Chris Tchou, Haarm-Pieter Duiker, Westley Sarokin, and Mark
322 Sagar. Acquiring the reflectance field of a human face. In *Proceedings of the 27th Annual*
323 *Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '00, page 145–156,
324 USA, 2000. ACM Press/Addison-Wesley Publishing Co.
- 325 [13] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt,
326 Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe
327 of annotated 3d objects. *arXiv preprint arXiv:2212.08051*, 2022.
- 328 [14] Zheng Ding, Xuaner Zhang, Zhihao Xia, Lars Jebe, Zhuowen Tu, and Xiuming Zhang. Dif-
329 fusionrig: Learning personalized priors for facial appearance editing. In *Proceedings of the*
330 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12736–12746, 2023.
- 331 [15] Epic Games. Unreal engine. <https://www.unrealengine.com>.
- 332 [16] Xiao Fu, Wei Yin, Mu Hu, Kaixuan Wang, Yuexin Ma, Ping Tan, Shaojie Shen, Dahua Lin, and
333 Xiaoxiao Long. Geowizard: Unleashing the diffusion priors for 3d geometry estimation from a
334 single image. *arXiv preprint arXiv:2403.12013*, 2024.
- 335 [17] David Futschik, Kelvin Ritland, James Vecore, Sean Fanello, Sergio Orts-Escolano, Brian
336 Curless, Daniel Šykora, and Rohit Pandey. Controllable light diffusion for portraits. In
337 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
338 8412–8421, 2023.
- 339 [18] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.
340 Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances*
341 *in Neural Information Processing Systems*, pages 6626–6637, 2017.
- 342 [19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances*
343 *in neural information processing systems*, 33:6840–6851, 2020.
- 344 [20] Michael Janner, Jiajun Wu, Tejas D Kulkarni, Ilker Yildirim, and Josh Tenenbaum. Self-
345 supervised intrinsic image decomposition. *Advances in neural information processing systems*,
346 30, 2017.
- 347 [21] Haian Jin, Isabella Liu, Peijia Xu, Xiaoshuai Zhang, Songfang Han, Sai Bi, Xiaowei Zhou,
348 Zexiang Xu, and Hao Su. Tensorir: Tensorial inverse rendering. In *Proceedings of the IEEE/CVF*
349 *Conference on Computer Vision and Pattern Recognition*, pages 165–174, 2023.
- 350 [22] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila.
351 Training generative adversarial networks with limited data. In *Proc. NeurIPS*, 2020.
- 352 [23] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative
353 adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and*
354 *pattern recognition*, pages 4401–4410, 2019.

- 355 [24] Bingxin Ke, Anton Obukhov, Shengyu Huang, Nando Metzger, Rodrigo Caye Daudt, and Kon-
356 rad Schindler. Repurposing diffusion-based image generators for monocular depth estimation.
357 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*
358 *(CVPR)*, 2024.
- 359 [25] Hoon Kim, Minje Jang, Wonjun Yoon, Jisoo Lee, Donghyun Na, and Sanghyun Woo. Switch-
360 light: Co-design of physics-driven architecture and pre-training framework for human portrait
361 relighting. *arXiv preprint arXiv:2402.18848*, 2024.
- 362 [26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
363 *arXiv:1412.6980*, 2014.
- 364 [27] Peter Kocsis, Julien Philip, Kalyan Sunkavalli, Matthias Nießner, and Yannick Hold-
365 Geoffroy. Lightit: Illumination modeling and control for diffusion models. *arXiv preprint*
366 *arXiv:2403.10615*, 2024.
- 367 [28] Peter Kocsis, Vincent Sitzmann, and Matthias Nießner. Intrinsic image diffusion for single-view
368 material estimation. *arXiv preprint arXiv:2312.12274*, 2023.
- 369 [29] Zhengfei Kuang, Kyle Olszewski, Menglei Chai, Zeng Huang, Panos Achlioptas, and Sergey
370 Tulyakov. Neroic: Neural rendering of objects from online image collections. *ACM Transactions*
371 *on Graphics (TOG)*, 41(4):1–12, 2022.
- 372 [30] Tianye Li, Timo Bolkart, Michael J. Black, Hao Li, and Javier Romero. Learning a model of
373 facial shape and expression from 4d scans. *ACM Trans. Graph.*, 36(6), nov 2017.
- 374 [31] Zhengqi Li and Noah Snavely. Learning intrinsic image decomposition from watching the
375 world. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
376 pages 9039–9048, 2018.
- 377 [32] Zhengqi Li, Zexiang Xu, Ravi Ramamoorthi, Kalyan Sunkavalli, and Manmohan Chandraker.
378 Learning to reconstruct shape and spatially-varying reflectance from a single image. *ACM Trans.*
379 *Graph.*, 37(6), dec 2018.
- 380 [33] Xian Liu, Jian Ren, Aliaksandr Siarohin, Ivan Skorokhodov, Yanyu Li, Dahua Lin, Xihui Liu,
381 Ziwei Liu, and Sergey Tulyakov. Hyperhuman: Hyper-realistic human generation with latent
382 structural diffusion. *arXiv preprint arXiv:2310.08579*, 2023.
- 383 [34] Yuan Liu, Peng Wang, Cheng Lin, Xiaoxiao Long, Jiepeng Wang, Lingjie Liu, Taku Komura,
384 and Wenping Wang. Nero: Neural geometry and brdf reconstruction of reflective objects from
385 multiview images. *ACM Transactions on Graphics (TOG)*, 42(4):1–22, 2023.
- 386 [35] Yunfei Liu, Yu Li, Shaodi You, and Feng Lu. Unsupervised learning for intrinsic image
387 decomposition from a single image. In *Proceedings of the IEEE/CVF conference on computer*
388 *vision and pattern recognition*, pages 3248–3257, 2020.
- 389 [36] Xiaoxiao Long, Yuan-Chen Guo, Cheng Lin, Yuan Liu, Zhiyang Dou, Lingjie Liu, Yuexin Ma,
390 Song-Hai Zhang, Marc Habermann, Christian Theobalt, et al. Wonder3d: Single image to 3d
391 using cross-domain diffusion. *arXiv preprint arXiv:2310.15008*, 2023.
- 392 [37] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoor-
393 thi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis.
394 *Communications of the ACM*, 65(1):99–106, 2021.
- 395 [38] Jacob Munkberg, Jon Hasselgren, Tianchang Shen, Jun Gao, Wenzheng Chen, Alex Evans,
396 Thomas Müller, and Sanja Fidler. Extracting triangular 3d models, materials, and lighting
397 from images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
398 *Recognition*, pages 8280–8290, 2022.
- 399 [39] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,
400 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning
401 robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.

- 402 [40] Rohit Pandey, Sergio Orts-Escolano, Chloe Legendre, Christian Haene, Sofien Bouaziz,
403 Christoph Rhemann, Paul E Debevec, and Sean Ryan Fanello. Total relighting: learning
404 to relight portraits for background replacement. *ACM Trans. Graph.*, 40(4):43–1, 2021.
- 405 [41] Bui Tuong Phong. Illumination for computer generated pictures. *Commun. ACM*, 18(6):311–317,
406 jun 1975.
- 407 [42] Puntawat Ponglertnapakorn, Nontawat Tritrong, and Supasorn Suwajanakorn. Difareli: Diffu-
408 sion face relighting. In *Proceedings of the IEEE/CVF International Conference on Computer*
409 *Vision*, pages 22646–22657, 2023.
- 410 [43] Mallikarjun B R, Ayush Tewari, Abdallah Dib, Tim Weyrich, Bernd Bickel, Hans-Peter Seidel,
411 Hanspeter Pfister, Wojciech Matusik, Louis Chevallier, Mohamed Elgharib, and Christian
412 Theobalt. Photoapp: photorealistic appearance editing of head portraits. *ACM Trans. Graph.*,
413 40(4), jul 2021.
- 414 [44] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
415 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
416 models from natural language supervision. In *International conference on machine learning*,
417 pages 8748–8763. PMLR, 2021.
- 418 [45] Mike Roberts, Jason Ramapuram, Anurag Ranjan, Atulit Kumar, Miguel Angel Bautista, Nathan
419 Paczan, Russ Webb, and Joshua M. Susskind. Hypersim: A photorealistic synthetic dataset for
420 holistic indoor scene understanding. In *International Conference on Computer Vision (ICCV)*
421 *2021*, 2021.
- 422 [46] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
423 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF*
424 *conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- 425 [47] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
426 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF*
427 *conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- 428 [48] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen.
429 Improved techniques for training gans. In *Advances in Neural Information Processing Systems*,
430 pages 2234–2242, 2016.
- 431 [49] Yizhi Song, Zhifei Zhang, Zhe Lin, Scott Cohen, Brian Price, Jianming Zhang, Soo Ye Kim,
432 and Daniel Aliaga. Objectstitch: Object compositing with diffusion model. In *Proceedings of*
433 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18310–18319,
434 2023.
- 435 [50] Pratul P Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and
436 Jonathan T Barron. Nerv: Neural reflectance and visibility fields for relighting and view
437 synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
438 *Recognition*, pages 7495–7504, 2021.
- 439 [51] Tiancheng Sun, Jonathan T Barron, Yun-Ta Tsai, Zexiang Xu, Xueming Yu, Graham Fyffe,
440 Christoph Rhemann, Jay Busch, Paul Debevec, and Ravi Ramamoorthi. Single image portrait
441 relighting. *ACM Transactions on Graphics (TOG)*, 38(4):1–12, 2019.
- 442 [52] Unity Technologies. Unity. <https://www.unity.com>.
- 443 [53] Bruce Walter, Stephen R. Marschner, Hongsong Li, and Kenneth E. Torrance. Microfacet models
444 for refraction through rough surfaces. In *Proceedings of the 18th Eurographics Conference on*
445 *Rendering Techniques*, EGSR’07, page 195–206, Goslar, DEU, 2007. Eurographics Association.
- 446 [54] Yifan Wang, Aleksander Holynski, Xiuming Zhang, and Xuaner Zhang. Sunstage: Portrait
447 reconstruction and relighting using the sun as a light stage. In *Proceedings of the IEEE/CVF*
448 *Conference on Computer Vision and Pattern Recognition*, pages 20792–20802, 2023.

- 449 [55] Zhibo Wang, Xin Yu, Ming Lu, Quan Wang, Chen Qian, and Feng Xu. Single image portrait
450 relighting via explicit multiple reflectance channel modeling. *ACM Transactions on Graphics*
451 (*TOG*), 39(6):1–13, 2020.
- 452 [56] Felix Wimbauer, Shangzhe Wu, and Christian Rupprecht. De-rendering 3d objects in the wild.
453 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
454 pages 18490–18499, 2022.
- 455 [57] Liwen Wu, Rui Zhu, Mustafa B Yaldiz, Yinhao Zhu, Hong Cai, Janarбек Matai, Fatih Porikli,
456 Tzu-Mao Li, Manmohan Chandraker, and Ravi Ramamoorthi. Factorized inverse path tracing
457 for efficient and accurate material-lighting estimation. In *Proceedings of the IEEE/CVF*
458 *International Conference on Computer Vision*, pages 3848–3858, 2023.
- 459 [58] Shangzhe Wu, Ameesh Makadia, Jiajun Wu, Noah Snavely, Richard Tucker, and Angjoo
460 Kanazawa. De-rendering the world’s revolutionary artefacts. In *Proceedings of the IEEE/CVF*
461 *conference on computer vision and pattern recognition*, pages 6338–6347, 2021.
- 462 [59] Shangzhe Wu, Christian Rupprecht, and Andrea Vedaldi. Unsupervised learning of probably
463 symmetric deformable 3d objects from images in the wild. In *Proceedings of the IEEE/CVF*
464 *conference on computer vision and pattern recognition*, pages 1–10, 2020.
- 465 [60] Binxin Yang, Shuyang Gu, Bo Zhang, Ting Zhang, Xuejin Chen, Xiaoyan Sun, Dong Chen,
466 and Fang Wen. Paint by example: Exemplar-based image editing with diffusion models. In
467 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
468 18381–18391, 2023.
- 469 [61] Yu-Ying Yeh, Koki Nagano, Sameh Khamis, Jan Kautz, Ming-Yu Liu, and Ting-Chun Wang.
470 Learning to relight portrait images via a virtual light stage and synthetic-to-real adaptation.
471 *ACM Transactions on Graphics (TOG)*, 41(6):1–21, 2022.
- 472 [62] Ye Yu and William AP Smith. Inverserendernet: Learning single image inverse rendering. In
473 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
474 3155–3164, 2019.
- 475 [63] Chong Zeng, Yue Dong, Pieter Peers, Youkang Kong, Hongzhi Wu, and Xin Tong. Di-
476 lightnet: Fine-grained lighting control for diffusion-based image generation. *arXiv preprint*
477 *arXiv:2402.11929*, 2024.
- 478 [64] Jingyang Zhang, Yao Yao, Shiwei Li, Jingbo Liu, Tian Fang, David McKinnon, Yanghai Tsin,
479 and Long Quan. Neilf++: Inter-reflectable light fields for geometry and material estimation. In
480 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3601–3610,
481 2023.
- 482 [65] Kai Zhang, Fujun Luan, Zhengqi Li, and Noah Snavely. Iron: Inverse rendering by optimizing
483 neural sdfs and materials from photometric images. In *Proceedings of the IEEE/CVF conference*
484 *on computer vision and pattern recognition*, pages 5565–5574, 2022.
- 485 [66] Kai Zhang, Fujun Luan, Qianqian Wang, Kavita Bala, and Noah Snavely. Physg: Inverse rendering
486 with spherical gaussians for physics-based material editing and relighting. In *Proceedings*
487 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5453–5462,
488 2021.
- 489 [67] Longwen Zhang, Qixuan Zhang, Minye Wu, Jingyi Yu, and Lan Xu. Neural video portrait
490 relighting in real-time via consistency modeling. In *Proceedings of the IEEE/CVF international*
491 *conference on computer vision*, pages 802–812, 2021.
- 492 [68] Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul Debevec, William T Freeman, and
493 Jonathan T Barron. Nerfactor: Neural factorization of shape and reflectance under an unknown
494 illumination. *ACM Transactions on Graphics (ToG)*, 40(6):1–18, 2021.
- 495 [69] Yuanqing Zhang, Jiaming Sun, Xingyi He, Huan Fu, Rongfei Jia, and Xiaowei Zhou. Modeling
496 indirect illumination for inverse rendering. In *Proceedings of the IEEE/CVF Conference on*
497 *Computer Vision and Pattern Recognition*, pages 18643–18652, 2022.

- 498 [70] Hao Zhou, Sunil Hadap, Kalyan Sunkavalli, and David W Jacobs. Deep single-image portrait
499 relighting. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages
500 7194–7202, 2019.
- 501 [71] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image
502 translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international
503 conference on computer vision*, pages 2223–2232, 2017.

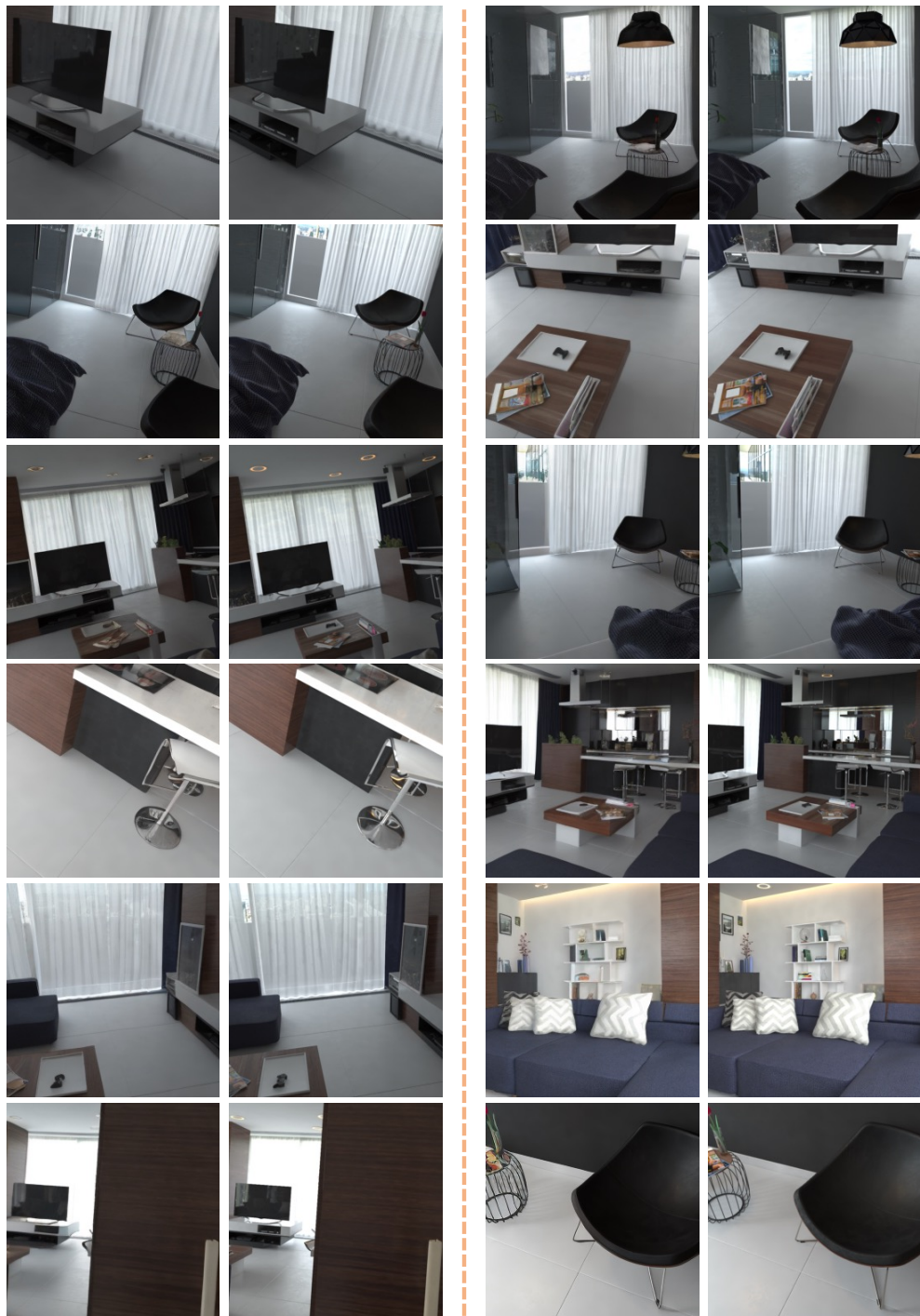


Figure 7: **Rendering results on Hypersim dataset.** The first and third columns are the rendered images while the second and the fourth are the ground truth images.

505 **NeurIPS Paper Checklist**

506 **1. Claims**

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

509 Answer: [Yes]

510 Justification: The main claims made in the abstract and introduction accurately reflect the
511 contributions and scope of our paper.

512 Guidelines:

- 513 • The answer NA means that the abstract and introduction do not include the claims
514 made in the paper.
- 515 • The abstract and/or introduction should clearly state the claims made, including the
516 contributions made in the paper and important assumptions and limitations. A No or
517 NA answer to this question will not be perceived well by the reviewers.
- 518 • The claims made should match theoretical and experimental results, and reflect how
519 much the results can be expected to generalize to other settings.
- 520 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
521 are not attained by the paper.

522 **2. Limitations**

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

524 Answer: [Yes]

525 Justification: Our paper discusses the limitations of our method.

526 Guidelines:

- 527 • The answer NA means that the paper has no limitation while the answer No means that
528 the paper has limitations, but those are not discussed in the paper.
- 529 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 530 • The paper should point out any strong assumptions and how robust the results are to
531 violations of these assumptions (e.g., independence assumptions, noiseless settings,
532 model well-specification, asymptotic approximations only holding locally). The authors
533 should reflect on how these assumptions might be violated in practice and what the
534 implications would be.
- 535 • The authors should reflect on the scope of the claims made, e.g., if the approach was
536 only tested on a few datasets or with a few runs. In general, empirical results often
537 depend on implicit assumptions, which should be articulated.
- 538 • The authors should reflect on the factors that influence the performance of the approach.
539 For example, a facial recognition algorithm may perform poorly when image resolution
540 is low or images are taken in low lighting. Or a speech-to-text system might not be
541 used reliably to provide closed captions for online lectures because it fails to handle
542 technical jargon.
- 543 • The authors should discuss the computational efficiency of the proposed algorithms
544 and how they scale with dataset size.
- 545 • If applicable, the authors should discuss possible limitations of their approach to
546 address problems of privacy and fairness.
- 547 • While the authors might fear that complete honesty about limitations might be used by
548 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
549 limitations that aren't acknowledged in the paper. The authors should use their best
550 judgment and recognize that individual actions in favor of transparency play an impor-
551 tant role in developing norms that preserve the integrity of the community. Reviewers
552 will be specifically instructed to not penalize honesty concerning limitations.

553 **3. Theory Assumptions and Proofs**

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

556 Answer: [NA]

557 Justification: Our paper does not include theoretical results.

558 Guidelines:

- 559 • The answer NA means that the paper does not include theoretical results.
- 560 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 561 referenced.
- 562 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 563 • The proofs can either appear in the main paper or the supplemental material, but if
- 564 they appear in the supplemental material, the authors are encouraged to provide a short
- 565 proof sketch to provide intuition.
- 566 • Inversely, any informal proof provided in the core of the paper should be complemented
- 567 by formal proofs provided in appendix or supplemental material.
- 568 • Theorems and Lemmas that the proof relies upon should be properly referenced.

569 4. Experimental Result Reproducibility

570 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

571 perimental results of the paper to the extent that it affects the main claims and/or conclusions

572 of the paper (regardless of whether the code and data are provided or not)?

573 Answer: [Yes]

574 Justification: Our paper fully discloses all the information needed to reproduce the main

575 experimental results.

576 Guidelines:

- 577 • The answer NA means that the paper does not include experiments.
- 578 • If the paper includes experiments, a No answer to this question will not be perceived
- 579 well by the reviewers: Making the paper reproducible is important, regardless of
- 580 whether the code and data are provided or not.
- 581 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 582 to make their results reproducible or verifiable.
- 583 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 584 For example, if the contribution is a novel architecture, describing the architecture fully
- 585 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 586 be necessary to either make it possible for others to replicate the model with the same
- 587 dataset, or provide access to the model. In general, releasing code and data is often
- 588 one good way to accomplish this, but reproducibility can also be provided via detailed
- 589 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 590 of a large language model), releasing of a model checkpoint, or other means that are
- 591 appropriate to the research performed.
- 592 • While NeurIPS does not require releasing code, the conference does require all submis-
- 593 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 594 nature of the contribution. For example
- 595 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
- 596 to reproduce that algorithm.
- 597 (b) If the contribution is primarily a new model architecture, the paper should describe
- 598 the architecture clearly and fully.
- 599 (c) If the contribution is a new model (e.g., a large language model), then there should
- 600 either be a way to access this model for reproducing the results or a way to reproduce
- 601 the model (e.g., with an open-source dataset or instructions for how to construct
- 602 the dataset).
- 603 (d) We recognize that reproducibility may be tricky in some cases, in which case
- 604 authors are welcome to describe the particular way they provide for reproducibility.
- 605 In the case of closed-source models, it may be that access to the model is limited in
- 606 some way (e.g., to registered users), but it should be possible for other researchers
- 607 to have some path to reproducing or verifying the results.

608 5. Open access to data and code

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

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

611 material?

612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663

Answer: [Yes]

Justification: We will release our code and data soon.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 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).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

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

Answer: [Yes]

Justification: All the details mentioned above are presented in the "Experiments" section.

Guidelines:

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

7. Experiment Statistical Significance

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

Answer: [NA]

Justification: Our paper does not include such experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- 664 • It is OK to report 1-sigma error bars, but one should state it. The authors should
665 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
666 of Normality of errors is not verified.
- 667 • For asymmetric distributions, the authors should be careful not to show in tables or
668 figures symmetric error bars that would yield results that are out of range (e.g. negative
669 error rates).
- 670 • If error bars are reported in tables or plots, The authors should explain in the text how
671 they were calculated and reference the corresponding figures or tables in the text.

672 8. Experiments Compute Resources

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

676 Answer: [Yes]

677 Justification: This information is provided in the "Experiments" section.

678 Guidelines:

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

687 9. Code Of Ethics

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

690 Answer: [Yes]

691 Justification: The research conducted in this paper conform, in every respect, with the
692 NeurIPS Code of Ethics.

693 Guidelines:

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

699 10. Broader Impacts

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

702 Answer: [NA]

703 Justification: Our paper does not discuss potential negative societal impacts of the work
704 performed.

705 Guidelines:

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

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

728 11. Safeguards

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

732 Answer: [NA]

733 Justification: Our paper poses no such risks.

734 Guidelines:

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

745 12. Licenses for existing assets

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

749 Answer: [Yes]

750 Justification: We cite the original papers and the website that produce the training and
751 evaluation dataset.

752 Guidelines:

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

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

768 13. **New Assets**

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

771 Answer: [Yes]

772 Justification: We introduce new assets in the "Experiments" section.

773 Guidelines:

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

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

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

786 Answer: [NA]

787 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

788 Guidelines:

- 789 • The answer NA means that the paper does not involve crowdsourcing nor research with
790 human subjects.
- 791 • Including this information in the supplemental material is fine, but if the main contribu-
792 tion of the paper involves human subjects, then as much detail as possible should be
793 included in the main paper.
- 794 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
795 or other labor should be paid at least the minimum wage in the country of the data
796 collector.

797 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 798 Subjects**

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

803 Answer: [NA]

804 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

805 Guidelines:

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