UrbanIR: Large-Scale Urban Scene Inverse Rendering from a Single Video Supplementary Material

Anonymous CVPR submission

Paper ID

1. Method

036

UrbanIR takes a multi-frame video of a scene under sin-002 gle illumination; as the camera moves, its motion is known. 003 Write $\{I_i, E_i, K_i\}$, where $I_i \in \mathbb{R}^{H \times W \times 3}$ is the RGB im-004 age; $E_i \in SE(3)$ is the camera pose; and K_i is camera 005 006 intrinsic matrix. We produce a neural field model that can be viewed from novel camera viewpoints under novel light-007 800 ing conditions. We do so by constructing a neural scene 009 model that encodes albedo, normal, semantics, and visibility in a unified manner (Sec. 1.1). This model is rendered 010 from a given camera pose with given illumination using an 011 012 end-to-end differentiable volume renderer (Sec. 1.2). Our 013 inference is by joint optimization of all properties (Sec. 1.3). Applications include changing the sun angle (??; top right), 014 015 day-to-night transitions (??; bottom right), and object in-016 sertion (??; middle right). More details about applications 017 are in Sec. 2. Fig. ?? provides an overview of our proposed inverse graphics and simulation framework. 018

019 1.1. Relightable Neural Scene Model

The scene representation is built on Instant-NGP [53, 58], 020 021 a spatial hash-based voxel NeRF representation. Instant-022 NGP offers numerous advantages, including low memory 023 consumption; high efficiency in training and rendering; and 024 compatibility with expansive outdoor scenes. Write $\mathbf{x} \in \mathbb{R}^3$ for position in 3D, d for query ray direction, θ for learnable 025 scene parameters; NeRF models, including Instant-NGP, 026 learn a radiance field $F_{\theta}(\mathbf{x}, \mathbf{d}) = (\mathbf{c}, \sigma)$, where $\mathbf{c} \in \mathbb{R}^3$ and 027 $\sigma \in \mathbb{R}$ represent observed color and opacity respectively. 028 029 Standard NeRFs have view- and lighting-dependent effects, such as shading, shadow, and specularity, baked into their 030 observed color, making them non-relightable. 031

032In contrast, UrbanIR learns a model of the intrinsic scene033attributes field independent of viewing angles and lighting034conditions. Write diffuse albedo a, surface normal n, seman-035tic vector s, and density σ ; then UrbanIR learns:

$$F_{\theta}(\mathbf{x}) = (\mathbf{a}, \mathbf{n}, \mathbf{s}, \sigma) \tag{1}$$

037 where θ is learnable parameters. The diffuse albedo rep-

resents the intrinsic color and texture of the material; the 038 normal represents the intrinsic surface geometry; density 039 encodes the spatial opacity, and semantics is used as a key to 040 query surface reflectance. Following Instant-NGP [53], we 041 learn a dense feature hash table to represent the scene, and 042 an individual MLP header is used to decode each attribute 043 given a queried feature at point x. We provide the details of 044 the architecture in the supplementary. The geometry of the 045 scene is implicitly encoded in σ . In contrast to existing re-046 lightable outdoor scene models that demand coupled explicit 047 geometry [60, 73], our scene model is implicit, providing 048 compactness and consistency to appearance modeling. 049

The lighting model is a parametric sun-sky model [34, 80].050This encodes outdoor illumination as:051

$$\mathbf{L} = \{ (\mathbf{L}_{sun}, \psi_{sun}, \phi_{sun}), \mathbf{L}_{amb}, \mathbf{L}_{sky} \}.$$
(2) 052

Our sun model is a 5-DoF representation, encoding sun 053 color \mathbf{L}_{sun} along with the azimuth and zenith ψ_{sun}, ϕ_{sun} . The 054 L_{amb} model is represented as a 3-DoF ambient light. The sky 055 dome model infers the sky texture from the viewing direction: 056 $C_{sky} = L_{sky}(\mathbf{r})$. We chose this minimalist sun-sky model as 057 it is more compact than other alternatives (e.g., HDR dome 058 or Spherical Gaussians) yet has proven highly effective in 059 modeling various outdoor illumination effects [34, 80]. 060

1.2. Rendering

Given the scene model F_{θ} and a lighting model **L**, rendering involves two steps: 1) volume rendering of the scene's intrinsic properties and visibility map onto the image plane, and 2) a shading process to produce the final result with view-dependent and lighting-dependent effects:

$$\mathbf{C} = \text{Shade}(\text{Intrinsic}(F_{\theta}, \mathbf{r}), \text{Shadow}(F_{\theta}, \mathbf{r}, \mathbf{L}), \mathbf{L})$$
(3)

where \mathbf{L} is the lighting model, \mathbf{C} is the final RGB color.

Intrinsics images are obtained by volume rendering. We accumulate predictions from $F(\cdot; \theta)$ along the query ray. 070 Multiple points are sampled along the ray, and intrinsics 071

061

062 063 064

065

066

067

at the query pixel along the ray [29, 51]. In particular, the albedo A, normal N, and semantics S are predicted as:

074
$$\mathbf{A}(\mathbf{r}) = \sum_{i=1}^{N} w_i \mathbf{a}_i, \mathbf{N}(\mathbf{r}) = \sum_{i=1}^{N} w_i \mathbf{n}_i, \mathbf{S}(\mathbf{r}) = \sum_{i=1}^{N} w_i \mathbf{s}_i, \quad (4)$$

where $w_i = \exp(-\sum_{j=1}^{i-1} \sigma_j \delta_j) (1 - \exp(-\sigma_i \delta_i))$ is alphacomposition weight, $\delta_i = t_i - t_{i-1}$. We perform rendering for each camera ray and get the final semantic map, albedo map, and the normal map.

Shadow modeling and rendering are essential for obtain-079 080 ing realistic-looking outdoor images. Modeling the visibility of the sun with a per-scene optimized MLP head (as 081 082 in [85, 87]) is impractical because we need to change the sun's position in relighting but can learn from only one po-083 sition. An alternative is to construct an explicit geometry 084 085 model to cast shadows [73], but this model might not be consistent with the other neural fields, and imposing consis-086 tency is difficult. Instead, we first compute an estimate $\mathbf{x}(\mathbf{r})$ 087 of the 3D point being shaded, then estimate the visibility 088 089 $V(\mathbf{x}, \mathrm{sun})$. Our key insight is that shadows in outdoor scenes 090 are primarily due to the visibility of a single directional sun-091 light.

We obtain $\mathbf{x}(\mathbf{r})$ for each ray by volume rendering depth (so substitute $\hat{t} = \sum w_i t_i$ into the equation for the ray being rendered). Now, to check whether \mathbf{x} is visible to the light source, we compute the transmittance along the ray segment between \mathbf{x} and the light source using volume rendering:

$$V(\mathbf{x}, \operatorname{sun}) = \exp\left(-\sum_{i} \sigma_{i}(\mathbf{x}_{i})\delta_{i}\right) \text{ where } \mathbf{x}_{i} = \mathbf{x} + t_{i}\mathbf{l}_{\operatorname{sun}}$$
(5)

098 Lower transmittance along a ray from a surface point 099 to a light source suggests fewer obstacles between the point and the light source. Eq. 5 establishes a strong link 100 between transmittance, lighting, and visibility fields used in 101 102 training. In particular, a point in a training image known 103 as shadowed (resp. out of shadow) should have large (resp. 104 small) accumulated transmittance. We use this constraint to adjust distant geometry during training. Compared to other 105 alternatives [73, 87], our proposed visibility test is simple to 106 compute, flexible for relighting, and aligns with intrinsic 107 108 properties with a few mild assumptions for outdoor scenes. 109

110Shading is performed by a Blinn-Phong model [7] that111incorporates sun and sky terms for the foreground scene and112an MLP query for the background sky. For $S(r) \in$ sky, we113use $C(r) = L_{sky}(r)$ and otherwise, we use

114
$$\mathbf{C}(\mathbf{r}) = \mathbf{A}(\mathbf{r}) \left(\mathbf{L}_{sun} \mathbf{D} \mathbf{V} + \mathbf{L}_{amb} \right)$$
(6)

where $\mathbf{D} = \max(\mathbf{N}(\mathbf{r}) \cdot \mathbf{l}_{sun}, 0)$ is the diffuse lighting at the surface, \mathbf{l}_{sun} is the sunlight direction (derived from ψ_{sun}, ϕ_{sun}). The visibility $V(\mathbf{x}, \operatorname{sun})$ is 1 if $\mathbf{x}(\mathbf{r})$ can see the sun and 0 otherwise. This shading model is capable of producing a realistic appearance with shadows following varying lighting conditions. The model can readily be extended with additional lighting sources at the relighting stage, as later shown in the night simulation. 117 118 119 120 120 121 122

1.3. Inverse graphics

We train scene $F(\cdot)$ (Eq. 1) and lighting L (Eq. 2) models 124 jointly using a loss: 125

$$\min_{\theta, \mathbf{L}} \mathcal{L}_{\text{render}} + \lambda_0 \mathcal{L}_{\text{deshadow}} + \lambda_1 \mathcal{L}_{\text{visibility}} + \lambda_2 \mathcal{L}_{\text{normal}} + \lambda_3 \mathcal{L}_{\text{semantics}} + \lambda_4 \mathcal{L}_{\text{reg}}$$
(7) 126

where individual loss terms are described below.

Rendering loss measures the agreement between observed128images and images rendered from the model using the training view and lighting, yielding $\mathcal{L}_{render} = \sum_{\mathbf{r}} ||\mathbf{C}_{gt}(\mathbf{r}) -$ 130 $\mathbf{C}(\mathbf{r})||_2^2$, where C is rendered color per ray, as defined in131Eq. 3, and \mathbf{C}_{gt} is the observed "ground-truth" color. Minimizing the rendering loss ensures our scene model can reproduce the observed images.134

Deshadowed rendering loss forces shadow effects out of 135 the estimated albedo. In particular, we compute a shadow-136 free version of an image using an off-the-shelf shadow detec-137 tion and removal network [13, 22] to obtain $C_{deshadow}$. We 138 then render that image from the model using the training 139 view and lighting, but assuming that every point can see the 140 sun (equivalently $V(\mathbf{x}, \operatorname{sun}) = 1$ for every \mathbf{x}). This yields 141 $\mathbf{C}'(\theta)$. We then measure the agreement between the two 142 to obtain $\mathcal{L}_{\text{deshadow}} = \sum_{\mathbf{r}} |\mathbf{C}_{\text{deshadow}} - \mathbf{C}'(\theta)|^2$. The combi-143 nation of this deshadowed rendering loss and the original 144 rendering loss directly gauges how the visibility map influ-145 ences rendering, and helps disentangle albedo and shadows. 146

Visibility loss exploits shadow detection to improve geom-147 etry estimates. A pixel that is known to be in shadow must 148 be at a point that cannot see the sun, so constraining geom-149 etry along a ray from that pixel to the sun. This loss could 150 be computed by simply comparing visibility V(sun) with 151 the shadow masks used for $\mathcal{L}_{deshadow}$. However, there are 152 challenges: first, computing visibility requires another vol-153 ume rendering per sample point; second, back-propagation 154 through volume rendering, shading, and visibility computa-155 tion forms a long, non-linear gradient chain, and optimiza-156 tion becomes difficult. Instead, we construct an intermediate 157 "guidance" visibility estimate $\hat{V}(\mathbf{r})$ which is an MLP head 158 trained to reproduce the shadow masks, and compute 159

$$\mathcal{L}_{\text{visibility}} = \sum_{\mathbf{r} \in \mathcal{R}} \text{CE}\left(M(\mathbf{r}), \hat{V}(\mathbf{r})\right) + \text{CE}\left(V(\mathbf{r}), \hat{V}(\mathbf{r})\right),$$
(8)

where $M(\mathbf{r})$ is the shadow mask at pixel \mathbf{r} , and CE(.,.) is a cross-entropy loss. Here the first term forces the (relatively easily trained) \hat{V} to agree with the shadow masks, and the second forces V to agree with \hat{V} . 164

126 127

160

199

200

201

202

203

204

205

206

207



Figure 1. **Night-time rendering.** In this sequence of images, the scene changes from daytime to night-time by introducing new light sources such as a headlight on a car and a street lamp. The top three and bottom three rows are from the same driving video, but at different times. Our decomposition method successfully removes dark shadows with sharp boundaries, resulting in a more realistic rendering of new light sources (such as streetlights and headlights) during night-time simulations. Our method is superior to Instruct-NeRF2NeRF [23], a data-driven, generative prior, radiance field approach.

165 166 167

168

169

170

171

172

185

186

Normal loss is computed by comparing results N_{gt} from an off-the-shelf normal estimator [18, 30] to the output of the normal MLP. Recall the camera is known for training scenes and write **r** for the pixel corresponding to 3D point $\mathbf{x}(\mathbf{r})$. An alternate estimate of the normal follows from the density field: $\hat{N}(\mathbf{r}) = -\frac{\nabla \sigma(\mathbf{x})}{\|\nabla \sigma(\mathbf{x})\|}$. Then our normal loss is given by:

$$\mathcal{L}_{\text{normal}} = \sum_{\mathbf{r} \in \mathcal{R}} \left(\|N_{\text{gt}}(\mathbf{r}) - N(\mathbf{r})\|^2 + \|N(\mathbf{r}) - \hat{N}(\mathbf{r})\|^2 \right).$$
(9)

We adopt smooth normal regularization from Ref-NeRF [69],
producing better density. The normal loss requires secondorder derivatives of the density during back-propagation. An
efficient implementation using the Hessian vector products
(HVP) enables memory-efficient computation.

178 Semantic loss is computed by comparing predicted seman-179 tics s with labels in the dataset [42]. We use an additional 180 loss to encourage high-depth values in the sky region, yield-181 ing: $\mathcal{L}_{semantics} = \sum CE(S_{gt}(\mathbf{r}), S(\mathbf{r})) - \sum D(\mathbf{r}).$

$$r \in \mathcal{R}$$
 $r \in sky$ 182A regularization term is used to regularize the albedo of183the scene and ambient light intensity. This is necessary due184to the ill-posed nature of our optimization process. However,

removing the hard shadow from the sunlight in the albedo

field A remains a challenge, particularly in urban driving

sequences. To address this challenge, we introduce a prior that ensures the ground albedo is homogeneous. This is important because the ground region typically shares a similar albedo value. More specifically, we first compute the average ground albedo $\bar{\mathbf{A}}_{g}$ from albedo \mathbf{A} and semantic S_{gt} and regularize the albedo using $\mathcal{L}_{albedo} = \sum_{\mathbf{r} \in \text{ground}} \|\mathbf{A}(\mathbf{r}) - \bar{\mathbf{A}}_{g}\|_{2}$. 192

We also calculate an *ambient regularization* term as $\|\mathbf{L}_{amb}\|_2$. We regularize the intensity of ambient light to avoid unnatural color shifts in the recovered albedo caused by a large intensity of ambient light. Our regularization term is thus $\mathcal{L}_{reg} = \mathcal{L}_{albedo} + \|\mathbf{L}_{amb}\|_2$. 197

2. Application Details

As intrinsics are recovered, UrbanIR can be rendered using any preferred source model. Natural uses are rendering scenes with different sun configurations and simulating night-time.

Simulating night-time proceeds by defining headlights and street lights, then illuminating with scene model considering specularity and lens flare. For sky regions $\mathbf{S}(\mathbf{r}) \in \text{sky}$, we use $\mathbf{C}(\mathbf{r}) = \mathbf{L}_{\text{sky}}(\mathbf{r})$ and otherwise, we use

$$\mathbf{A}(\mathbf{r})\left(\sum \mathbf{L}_{dif}^{i}\mathbf{D}_{i}\mathbf{V}_{i}+\mathbf{L}_{amb}\right)+\sum_{i}\mathbf{L}_{spec}^{i}$$
(10) 208

288

209 210

212

229

248

The spotlight we used is given by the center $\mathbf{o}_L^i \in \mathbb{R}^3$ and direction $\mathbf{d}_{L}^{i} \in \mathbb{R}^{3}$ of the light. This spotlight produces a diffuse radiance at r given by 211

$$\mathbf{L}_{\text{dif}}^{i}(\mathbf{r}) = \frac{1}{\|\mathbf{o}_{L}^{i} - \mathbf{x}(\mathbf{r})\|^{2}} \left(l \cdot \mathbf{d}_{L}^{i}\right)^{k}, l = \frac{\mathbf{o}_{L}^{i} - \mathbf{x}(\mathbf{r})}{\|\mathbf{o}_{L}^{i} - \mathbf{x}(\mathbf{r})\|},$$
(11)

Spotlight's diffuse color intensity is brightest on the central 213 ray $\mathbf{r}(t) = \mathbf{o}_L - t\mathbf{d}_L$, decays with distance from ray $\mathbf{r}(t)$ 214 and angle. We modulate it with constant k. 215

The realistic night-time simulation requires reproducing 216 217 the strong specular effects on cars. We find car regions using a semantic field S in Eq. 4, then simulate specular reflec-218 219 tion with the Blinn-Phong model [7], where the γ (specular strength) parameter is inherited from the semantic field. 220

221 At night, luminaires often display lens flares. A pure 222 simulation of lens flares is impractical, as it requires extensive ray tracing through the lens. We use the standard 223 image-based approximation [1] to simulate such light 224 225 scattering effects. For directly visible luminaires, we composite a real-world lens flare image from a similar 226 lighting source into the image, using location and depth. As 227 228 Fig. 1 shows, this simple method is effective.

Object insertion proceeds by a hybrid rendering strategy. 230 We first cast rays from the camera and estimate ray-mesh 231 232 intersections [16] for the inserted object. If the ray hits the 233 mesh and the distance is shorter than the volume rendering depth, the albedo $A(\mathbf{r})$, normal $N(\mathbf{r})$, and depth $D(\mathbf{r})$ are 234 replaced with the object attributes. In the shadow pass, 235 we calculate visibility from surface points to the light 236 237 source (Eq. 5), and also estimate the ray-mesh intersection 238 for the tracing rays. If the rays hit the mesh (meaning occlusion by the object), the visibility is also updated 239 : $V(\mathbf{r}) = 0$. With updated $A(\mathbf{r}), N(\mathbf{r}), V(\mathbf{r})$, shading 240 (Eq. 1.2) is applied to render images with virtual objects. 241 242 Our method not only casts object shadows in the scene but 243 also casts scene shadows on the object, enhancing realism significantly. Similar approaches have been depicted in 244 recent works [37, 57]. However, ours is the first to be 245 visibility-aware, enabling us to render effects when an object 246 247 enters into a shadow.

Outdoor relighting is done by simply adjusting lighting pa-249 rameters (position or color of the sun; sky color) then re-250 251 rendering using Eq. 3. We also use semantics to interpret 252 specular car surfaces and emulate their reflectance during 253 the simulation.

3. Model Architecture 254

Instant-NGP [53] encodes the scene with a multi-scale hash 255 table, and each entry contains learnable parameters. For 256 257 point $\mathbf{x} \in \mathbb{R}^3$, the model retrieves and interpolates the parameters with hash function: $F(\mathbf{x}, \theta)$. UrbanIR adopts the 258 hash encoding from [53] and maintain two separate hash 259 tables for geometry and appearance, and predict the scene 260 properties with: 261

$$\sigma = F_g(\mathbf{x}, \theta_g)$$

$$(\mathbf{a}, \mathbf{n}, s) = F_a(\mathbf{x}, \theta_a),$$
(12) 262

where σ is density, $(\mathbf{a}, \mathbf{n}, s)$ are albedo, surface normal, and 263 semantic. θ_q , θ_a are learnable parameters for geometry and 264 appearance. Please note that the density field σ is not only 265 involved in the volume rendering (Eq. 4), but also involved 266 in visibility estimation (Eq. 5) and normal loss calculation 267 (Eq. 9). The hash encoding is implemented with tiny-cuda-268 nn [52]. We empirically find that maintaining separate learn-269 able parameters for geometry and appearance leads to more 270 stable convergence and higher rendering quality. 271

4. Training Details

The training procedure is illustrated in Fig. 2. We leverage 273 pretrained networks as 2D priors during training to address 274 the ill-posed inverse problem. Specifically, the shadow mask 275 is estimated with MTMT [13], and shadow removal is per-276 formed with ShadowFormer [22]. Omnidata normal esti-277 mation [18] helps refine scene geometry (Eq. 9), which is 278 critical in the shading quality and albedo decomposition. A 279 semantic map is provided in Kitti360 dataset [42] and can 280 also be estimated with MMSegmentation [15] if such infor-281 mation is not provided. The loss terms are weighted during 282 training: 283

$$\mathcal{L} = \mathcal{L}_{render} + \lambda_0 \mathcal{L}_{deshadow} + \lambda_1 \mathcal{L}_{visibility} + \lambda_2 \mathcal{L}_{normal} + \lambda_3 \mathcal{L}_{semantics} + \lambda_4 \mathcal{L}_{reg}$$

where $\lambda_0 = 1.0, \lambda_1 = 0.001, \lambda_2 = 0.01, \lambda_3 = 0.04, \lambda_4 =$ 285 0.1. We use Adam optimizer [31] with a learning rate of 286 0.002 for a total of 100 epochs during the optimization. 287

5. Related Works Comparison

We compare the problem setting, input requirement with 289 recent methods in Tab. 1. UrbanIR addresses inverse 290 rendering for large-scale urban scenes that object-centric 291 methods [28, 85, 87] fails to reconstruct. Furthermore, our 292 method takes videos under single illuminations as input. In 293 order to estimate the geometry of large-scale driving scenes, 294 FEGR [73] and LightSim [56] rely on captures from five 295 to six cameras and LiDAR sensors. On the other hand, Ur-296 banIR only needs videos from single or stereo cameras 297 without any guidance from LiDAR. Our method also per-298 forms nighttime simulation by inserting local light sources 299 (e.g. streetlight, vehicle light), which is not demonstrated in 300 previous works. 301



Figure 2. **Training Pipeline.** UrbanIR retrieves scene intrinsics with volume rendering from camera rays, which is guided by semantic and normal priors. Transmittance along tracing rays is supervised with shadow masks. The shading model (illustrated in Fig. ??) is performed **with** and **without** visibility term, and enforce reconstruction loss with original and deshadowed images, respectively. Please refer to Section 1.3 for more details.

302 6. Visibility Modeling for Object Insertion

Following [72, 73], we build the object insertion pipeline 303 304 with Blender [14], and the results are shown in Fig. 3. By tracing the rays from object surface toward light sources 305 (i.e. the sun), UrbanIR estimates the visibility with volume 306 307 rendering (Eq. 5 in main paper). As a result, our full model is able to cast scene shadow on the inserted objects and also 308 weaken the object shadow on the ground if it overlaps with 309 310 the existing scene shadow. The visibility modeling enables 311 our method to simulate shadows better and to enhance the 312 insertion realism significantly.

313 7. More Relighting Results

We compare the relighting quality with FEGR [73] in Fig. 4. 314 315 FEGR [73] first extracts mesh and estimates the shading 316 from the lighting configuration, and the imperfect mesh ge-317 ometry produces artifacts and loses appearance details. On 318 the other hand, our method alleviates the original shadow and produces relighting images while preserving appearance 319 320 details. Additionally, UrbanIR can insert local light sources 321 and simulate the scene in the night time (Fig. 5), demonstrat-322 ing the capability and flexibility of our relighting framework.

We show additional night simulation results on various Kitti360 [42] sequences in Fig. 6, demonstrating the generalization capability of UrbanIR . The Instruct-Pix2Pix [10] leverages the large language model [11] and stable diffusion [59] for abundant image editing tasks. However, such a data-driven method cannot move the daylight shading and 328 shadow in the input images. On the contrary, UrbanIR de-329 composes shadow-free albedo and performs physically-330 based rendering with new light sources (e.g., streetlights, 331 headlights), significantly enhancing the visual quality of 332 night simulation. The strong specular reflection is also sim-333 ulated on the car region (Eq. 10), boosting the realism of 334 metal material. Please note that the simulation is flexible, 335 and the user can adjust physical parameters (e.g., light color, 336 light strength) to create various effects. Please refer to our 337 supplementary video to visualize view consistency and con-338 trollable simulation better. 339

8. Baseline Details

Description of the approach of baselines we compared to. 341

NeRF + Mesh The recent work FEGR [73] explores the 342 relighting of outdoor scenes under singular or multiple il-343 lumination sources. However, due to the absence of open-344 source access to their method, we implement our baseline 345 model incorporating similar visibility modeling strategies. 346 Specifically, we employ the marching cubes technique [48] 347 to extract the mesh from our model, excluding our proposed 348 visibility optimization (as per Eq.8). In alignment with the 349 shadow mapping approach adopted by FEGR [73], we cast 350 shadows by estimating two intersections: the first between 351 the camera rays and the mesh and the second by tracing rays 352



Figure 3. **Object Insertion Qualitative Results.**Without visibility modeling (middle column), the scenes do not cast shadows on the inserted objects and the original object shadow looks unrealistic in the existing shadow. Our full method (right column) simulates the better interaction between the reconstructed scenes and inserted objects with the help of visibility modeling.

353 from the surface to the light source.

Colmap MVS [64] We compare our method with an
explicit geometry-based baseline. For this, we utilize
COLMAP for dense scene reconstruction [63, 64] and import
the resulting scene into Blender [14] for relighting simulation.

Instruct-Pix2Pix [10] edits images according to user instruction. The model leverages large language model GPT3 [11] and Stable Diffusion [59] for generating image and
instruction pairs and fine-tune diffusion model to perform
editing. We use instructions "change to night", and "It's now
midnight" for night image generation.

Instruct-NeRF2NeRF [23] aims to edit NeRF scenes
with text instructions. It uses a generative image editing
model [10] to iteratively edit input images while optimizing
the underlying scene model, resulting in an optimized 3D
scene that respects the instruction. We compare Instruct

NeRF2NeRF in night simulation, where we provide the instruction, "*Make it look like it was taken at night*." 371

NeRF-OSR [60] is a recent work for outdoor scene re-372 construction and relighting. We use the open-source project 373 provided by the author to run this baseline. This method 374 represents lighting as spherical harmonics parameters. It 375 is worth noting that NeRF-OSR was designed for inverse 376 rendering in multi-illumination conditions. For a fair com-377 parison, we rotate the spherical vectors to simulate different 378 light conditions. 379

RelightNet [79]is a single-image based relighting frame-work. We use the open-source project provided by the au-381thors to produce intrinsic decomposition results, including382shading and albedo for comparison.383



Figure 4. Relighting Comparison on Waymo Open Dataset [68]. The second and third columns compare the relighting quality. The authors provide the FEGR results and we match the lighting condition according to the shadow direction.



Figure 5. Controllable Night Simulation.

384 9. Related Work

Inverse Graphics involves inferring illumination and in-385 386 trinsic properties of a scene. The problem is underconstrained, and there is much reliance on priors [2, 3, 26, 387 27, 35, 50, 62, 78, 83] or on managed lighting condi-388 tions [2, 2, 20, 25, 25, 82], known geometry [17, 32, 36, 61], 389 or material simplifications [49, 83, 88]. Recent methods 390 use deep learning techniques to reason about material prop-391 392 erties [45-47, 55, 77, 86]. Models trained on synthetic

data [43] or pair-wise annotated data [4] have shown promis-393 ing results. Learned predictors of albedo or shading are de-394 scribed and reviewed in [6, 19, 65]. Neural representations 395 of material or illumination appear in [5, 38-41, 47]. Like 396 these methods, we exploit monocular cues, such as shadows 397 and surface normals. In contrast, we combine learning-based 398 monocular cues and model-based relightable NeRF optimiza-399 tion to infer the scene's intrinsic properties and illumination. 400



Figure 6. **Nighttime rendering.** The scene is transformed from daytime (1st row) to night-time (3rd row) by introducing new light sources: a headlight on a car and a street lamp. Top 3 and bottom 3 rows are from same driving sequence with different time stamp. Comparing with data-driven generative model and Instruct-Pix2Pix [10], the dark shadows with sharp boundaries are successfully removed with our decomposition, resulting more realistic rendering with new light sources (e.g. streetlights, headlight) during the nighttime simulation.

Method	Scene Type	Illumination Conditions	RGB Only	Nighttime Simulation
NeRFFactor [85]	Object	Multi	Yes	
TensoIR [28]	Object	Single	Yes	
InvRender [87]	Object	Single	Yes	
NeRF-OSR [60]	Front-Facing	Multi	Yes	
FEGR [73]	Large Scene	Single/Multi	LiDAR	
LightSim [56]	Large Scene	Single/Multi	LiDAR	
UrbanIR (Ours)	Large Scene	Single	Yes	\checkmark

Table 1. Comparison of various recent relightable NeRF methods. UrbanIR is among the first to offer single-illumination and RGBonly relightable NeRF capabilities suitable for large-scale scenes.

401 Relightable Neural Fields: Relightable neural radiance
402 field methods [8, 9, 24, 54, 74, 77, 81, 85] aim to factor the

neural field into multiple intrinsic components and lever-403 age neural shading equations for illumination and material 404 modeling. These methods admit realistic and controllable 405 rendering of scenes with varying lighting conditions and 406 materials. However, most relightable NeRF methods focus 407 on objects with surrounding views or small bounded indoor 408 environments. Important exceptions are: NeRF-OSR [60], 409 which assumes access to multiple lighting sources for de-410 composition, and NeRF meets explicit geometry [74], which 411 either uses multiple lighting or exploits depth sensing, such 412 as LiDAR. In contrast, our proposed approach only requires 413 a single video captured under the same, unknown illumina-414 tion, making it more applicable to a broader range of scenes. 415

Differentiable rendering techniques make inverse graph-

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

ics tasks more flexible and convenient [12, 44, 55]. Most
render meshes are suitable for object-level rendering, and so
are challenging to apply to large urban scenes. In contrast,
we leverage neural radiance fields (NeRF) [51].

Shadow modeling using images is challenging. Methods 421 422 trained to cast shadows from images [71, 84] are tailored for particular objects (pedestrians, cars, etc). Learned methods 423 424 can detect and remove shadows from 2D images [21, 22, 70]. But inverse graphics require modeling the full 3D geometry, 425 426 intrinsic scene properties, and ensuring temporal consistency. 427 Model-based optimization methods can infer shadows but 428 rely on accurate scene geometry [33, 67, 75]. Using visi-429 bility fields to model shadows results in difficulty provid-430 ing consistent shadows in relation to the underlying geom-431 etry [60, 66, 76, 87]. In contrast, our method combines the strengths of learning-based monocular shadow prediction 432 433 and removal and model-based inverse graphics.

434 References

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

- [1] Tomas Akenine-Moller, Eric Haines, and Naty Hoffman. *Real-time rendering*. AK Peters/crc Press, 2019. 4
- [2] Jonathan T Barron and Jitendra Malik. Shape, illumination, and reflectance from shading. *TPAMI*, 2014. 7
- [3] H.G. Barrow and Joan M. Tenenbaum. Recovering intrinsic scene characteristics from images. *Computer Vision Systems*, 1978. 7
- [4] Sean Bell, Kavita Bala, and Noah Snavely. Intrinsic images in the wild. 2014. 7
- [5] Anand Bhattad and D. A. Forsyth. Stylitgan: Prompting stylegan to produce new illumination conditions, 2023. 7
- [6] Anand Bhattad, Daniel McKee, Derek Hoiem, and DA Forsyth. Stylegan knows normal, depth, albedo, and more. arXiv preprint arXiv:2306.00987, 2023. 7
- [7] James F Blinn. Models of light reflection for computer synthesized pictures. In *Proceedings of the 4th annual conference* on Computer graphics and interactive techniques, pages 192– 198, 1977. 2, 4
- [8] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T. Barron, Ce Liu, and Hendrik P.A. Lensch. Nerd: Neural reflectance decomposition from image collections. In *ICCV*, 2021. 8
- [9] Mark Boss, Varun Jampani, Raphael Braun, Ce Liu,
 Jonathan T. Barron, and Hendrik P.A. Lensch. Neural-PIL:
 Neural Pre-Integrated Lighting for Reflectance Decomposition. In *Advances in Neural Information Processing Systems*(*NeurIPS*), 2021. 8
- 462 [10] Tim Brooks, Aleksander Holynski, and Alexei A Efros. In463 structpix2pix: Learning to follow image editing instructions.
 464 In *CVPR*, 2023. 5, 6, 8
- [11] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan,
 Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language
 models are few-shot learners. *NeurIPS*, 2020. 5, 6

- [12] Wenzheng Chen, Huan Ling, Jun Gao, Edward Smith, Jaakko
 Lehtinen, Alec Jacobson, and Sanja Fidler. Learning to predict 3d objects with an interpolation-based differentiable renderer. Advances in neural information processing systems, 32, 2019. 9
 473
- [13] Zhihao Chen, Lei Zhu, Liang Wan, Song Wang, Wei Feng, and Pheng-Ann Heng. A multi-task mean teacher for semisupervised shadow detection. In *CVPR*, 2020. 2, 4
- [14] Blender Online Community. Blender a 3D modelling and rendering package. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. 5, 6
- [15] MMSegmentation Contributors. MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark. https://github.com/openmmlab/mmsegmentation, 2020. 4
- [16] Dawson-Haggerty et al. trimesh. 4
- [17] Yue Dong, Guojun Chen, Pieter Peers, Jiawan Zhang, and Xin Tong. Appearance-from-motion: Recovering spatially varying surface reflectance under unknown lighting. ACM Transactions on Graphics (TOG), 33(6):1–12, 2014. 7
- [18] Ainaz Eftekhar, Alexander Sax, Jitendra Malik, and Amir Zamir. Omnidata: A scalable pipeline for making multi-task mid-level vision datasets from 3d scans. In *ICCV*, 2021. 3, 4
- [19] David Forsyth and Jason J Rock. Intrinsic image decomposition using paradigms. *IEEE transactions on pattern analysis and machine intelligence*, 44(11):7624–7637, 2021. 7
- [20] Roger Grosse, Micah K Johnson, Edward H Adelson, and William T Freeman. Ground truth dataset and baseline evaluations for intrinsic image algorithms. In *ICCV*, 2009. 7
- [21] Lanqing Guo, Chong Wang, Wenhan Yang, Siyu Huang, Yufei Wang, Hanspeter Pfister, and Bihan Wen. Shadowdiffusion: When degradation prior meets diffusion model for shadow removal. *arXiv preprint arXiv:2212.04711*, 2022. 9
- [22] Lanqing Guo, Siyu Huang, Ding Liu, Hao Cheng, and Bihan Wen. Shadowformer: Global context helps image shadow removal. *AAAI*, 2023. 2, 4, 9
- [23] Ayaan Haque, Matthew Tancik, Alexei Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-nerf2nerf: Editing 3d scenes with instructions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023. 3, 6
- [24] Jon Hasselgren, Nikolai Hofmann, and Jacob Munkberg.
 Shape, Light, and Material Decomposition from Images using Monte Carlo Rendering and Denoising. *arXiv:2206.03380*, 2022. 8
- [25] Daniel Hauagge, Scott Wehrwein, Kavita Bala, and Noah Snavely. Photometric ambient occlusion. In *CVPR*, 2013. 7
- [26] Berthold KP Horn. Determining lightness from an image. Computer graphics and image processing, 1974. 7
- [27] Berthold KP Horn. Obtaining shape from shading information. *The psychology of computer vision*, 1975. 7
- [28] Haian Jin, Isabella Liu, Peijia Xu, Xiaoshuai Zhang, Songfang Han, Sai Bi, Xiaowei Zhou, Zexiang Xu, and Hao Su. Tensoir: Tensorial inverse rendering. *CVPR*, 2023. 4, 8
- [29] James T Kajiya and Brian P Von Herzen. Ray tracing volume densities. *ACM SIGGRAPH computer graphics*, 1984. 2
- [30] Oğuzhan Fatih Kar, Teresa Yeo, Andrei Atanov, and Amir Zamir. 3d common corruptions and data augmentation. In *CVPR*, 2022. 3
 526

530

531

532

540

541

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

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

- 527 [31] Diederik P. Kingma and Jimmy Ba. Adam: A method for528 stochastic optimization. *ICLR*, 2014. 4
 - [32] Pierre-Yves Laffont, Adrien Bousseau, and George Drettakis. Rich intrinsic image decomposition of outdoor scenes from multiple views. *IEEE transactions on visualization and computer graphics*, 2012. 7
- [33] Samuli Laine, Timo Aila, Ulf Assarsson, Jaakko Lehtinen,
 and Tomas Akenine-Möller. Soft shadow volumes for ray
 tracing. In *ACM SIGGRAPH 2005 Papers*, pages 1156–1165.
 2005. 9
- [34] Jean-François Lalonde and Iain Matthews. Lighting estimation in outdoor image collections. In *International Conference on 3D Vision (3DV)*. IEEE, 2014. 1
 - [35] Edwin H Land and John J McCann. Lightness and retinex theory. Josa, 1971. 7
- [36] Hendrik PA Lensch, Jan Kautz, Michael Goesele, Wolfgang
 Heidrich, and Hans-Peter Seidel. Image-based reconstruction
 of spatial appearance and geometric detail. *TOG*, 2003. 7
- [37] Yuan Li, Zhi-Hao Lin, David Forsyth, Jia-Bin Huang, and
 Shenlong Wang. Climatenerf: Extreme weather synthesis
 in neural radiance field. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3227–
 3238, 2023. 4
- [38] Zhengqin Li, Kalyan Sunkavalli, and Manmohan Chandraker.
 Materials for masses: SVBRDF acquisition with a single
 mobile phone image. In *ECCV*, pages 72–87, 2018. 7
- [39] Zhengqin Li, Zexiang Xu, Ravi Ramamoorthi, Kalyan
 Sunkavalli, and Manmohan Chandraker. Learning to reconstruct shape and spatially-varying reflectance from a single
 image. ACM Transactions on Graphics (TOG), 37(6):1–11,
 2018.
- [40] Zhengqin Li, Mohammad Shafiei, Ravi Ramamoorthi, Kalyan
 Sunkavalli, and Manmohan Chandraker. Inverse rendering
 for complex indoor scenes: Shape, spatially-varying lighting
 and svbrdf from a single image. In *CVPR*, 2020.
- [41] Zhengqin Li, Jia Shi, Sai Bi, Rui Zhu, Kalyan Sunkavalli,
 Miloš Hašan, Zexiang Xu, Ravi Ramamoorthi, and Manmohan Chandraker. Physically-based editing of indoor scene
 lighting from a single image. ECCV, 2022. 7
- [42] Yiyi Liao, Jun Xie, and Andreas Geiger. KITTI-360: A novel dataset and benchmarks for urban scene understanding in 2d and 3d. *in arXiv*, 2021. 3, 4, 5
- [43] Daniel Lichy, Jiaye Wu, Soumyadip Sengupta, and David W
 Jacobs. Shape and material capture at home. In *Proceedings*of the IEEE/CVF Conference on Computer Vision and Pattern
 Recognition, pages 6123–6133, 2021. 7
- 573 [44] Shichen Liu, Tianye Li, Weikai Chen, and Hao Li. Soft raster574 izer: A differentiable renderer for image-based 3d reasoning.
 575 In Proceedings of the IEEE/CVF International Conference
 576 on Computer Vision, pages 7708–7717, 2019. 9
- 577 [45] Stephen Lombardi and Ko Nishino. Reflectance and illumi578 nation recovery in the wild. *IEEE transactions on pattern*579 *analysis and machine intelligence*, 38(1):129–141, 2015. 7
- [46] Stephen Lombardi and Ko Nishino. Radiometric scene decomposition: Scene reflectance, illumination, and geometry from rgb-d images. In 2016 Fourth International Conference on 3D Vision (3DV), pages 305–313. IEEE, 2016.

- [47] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. ACM Transactions on Graphics (TOG), 38(4):65, 2019. 7
- [48] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. ACM siggraph computer graphics, 1987. 5
- [49] Wei-Chiu Ma, Hang Chu, Bolei Zhou, Raquel Urtasun, and Antonio Torralba. Single image intrinsic decomposition without a single intrinsic image. In *ECCV*, 2018. 7
- [50] Stephen Robert Marschner. Inverse rendering for computer graphics. Cornell University, 1998. 7
- [51] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In ECCV, 2020. 2, 9
- [52] Thomas Müller. tiny-cuda-nn, 2021. 4
- [53] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM TOG*, 2022. **1**, 4
- [54] Jacob Munkberg, Jon Hasselgren, Tianchang Shen, Jun Gao, Wenzheng Chen, Alex Evans, Thomas Müller, and Sanja Fidler. Extracting Triangular 3D Models, Materials, and Lighting From Images. arXiv:2111.12503 [cs], 2021. 8
- [55] Jacob Munkberg, Jon Hasselgren, Tianchang Shen, Jun Gao, Wenzheng Chen, Alex Evans, Thomas Müller, and Sanja Fidler. Extracting triangular 3d models, materials, and lighting from images. In *CVPR*, 2022. 7, 9
- [56] Ava Pun, Gary Sun, Jingkang Wang, Yun Chen, Ze Yang, Sivabalan Manivasagam, Wei-Chiu Ma, and Raquel Urtasun. Neural lighting simulation for urban scenes. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
 4, 8
- [57] Yi-Ling Qiao, Alexander Gao, Yiran Xu, Yue Feng, Jia-Bin Huang, and Ming C Lin. Dynamic mesh-aware radiance fields. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 385–396, 2023. 4
- [58] Chen Quei-An. ngp_pl: a pytorch-lightning implementation of instant-ngp, 2022. 1
- [59] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In CVPR, 2022. 5, 6
- [60] Viktor Rudnev, Mohamed Elgharib, William Smith, Lingjie Liu, Vladislav Golyanik, and Christian Theobalt. Nerf for outdoor scene relighting. In *European Conference on Computer Vision (ECCV)*, 2022. 1, 6, 8, 9
- [61] Imari Sato, Yoichi Sato, and Katsushi Ikeuchi. Illumination from shadows. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2003. 7
- [62] Yoichi Sato, Mark D Wheeler, and Katsushi Ikeuchi. Object shape and reflectance modeling from observation. In *SIGGRAPH*, 1997. 7
- [63] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 6
- [64] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *ECCV*, 2016. 6

704

705

706

707

708

709

710

711

712

713

714

715

716

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

- [65] Soumyadip Sengupta, Jinwei Gu, Kihwan Kim, Guilin Liu,
 David W Jacobs, and Jan Kautz. Neural inverse rendering
 of an indoor scene from a single image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
 pages 8598–8607, 2019. 7
- [66] Pratul P Srinivasan, Boyang Deng, Xiuming Zhang, Matthew
 Tancik, Ben Mildenhall, and Jonathan T Barron. Nerv: Neural reflectance and visibility fields for relighting and view
 synthesis. In *CVPR*, 2021. 9
- [67] Jon Story. Hybrid ray traced shadows. In *Game Developer Conference*, 2015. 9
- [68] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien 652 653 Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, 654 Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, 655 Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott Et-656 tinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang, 657 Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. 658 Scalability in perception for autonomous driving: Waymo 659 open dataset. In CVPR, 2020. 7
- [69] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler,
 Jonathan T. Barron, and Pratul P. Srinivasan. Ref-NeRF:
 Structured view-dependent appearance for neural radiance
 fields. *CVPR*, 2022. 3
- [70] Jin Wan, Hui Yin, Zhenyao Wu, Xinyi Wu, Yanting Liu, and
 Song Wang. Style-guided shadow removal. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XIX*, pages
 361–378. Springer, 2022. 9
- [71] Yifan Wang, Andrew Liu, Richard Tucker, Jiajun Wu, Brian L
 Curless, Steven M Seitz, and Noah Snavely. Repopulating
 street scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5110–5119,
 2021. 9
- [72] Zian Wang, Wenzheng Chen, David Acuna, Jan Kautz, and
 Sanja Fidler. Neural light field estimation for street scenes
 with differentiable virtual object insertion. In *ECCV*, 2022. 5
- [73] Zian Wang, Tianchang Shen, Jun Gao, Shengyu Huang, Jacob
 Munkberg, Jon Hasselgren, Zan Gojcic, Wenzheng Chen, and
 Sanja Fidler. Neural fields meet explicit geometric representations for inverse rendering of urban scenes. In *CVPR*, 2023.
 1, 2, 4, 5, 8
- [74] Zian Wang, Tianchang Shen, Jun Gao, Shengyu Huang, Jacob
 Munkberg, Jon Hasselgren, Zan Gojcic, Wenzheng Chen, and
 Sanja Fidler. Neural fields meet explicit geometric representation for inverse rendering of urban scenes. *arXiv*, 2023.
 8
- [75] Tai-Pang Wu, Chi-Keung Tang, Michael S Brown, and Heung-Yeung Shum. Natural shadow matting. *ACM Transactions on Graphics (TOG)*, 26(2):8–es, 2007.
- [76] Wenqi Yang, Guanying Chen, Chaofeng Chen, Zhenfang
 Chen, and Kwan-Yee K Wong. S 3-nerf: Neural reflectance
 field from shading and shadow under a single viewpoint. *arXiv preprint arXiv:2210.08936*, 2022. 9
- [77] Ye Yu and William AP Smith. Inverserendernet: Learning single image inverse rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3155–3164, 2019. 7, 8

- [78] Yizhou Yu, Paul Debevec, Jitendra Malik, and Tim Hawkins.
 Inverse global illumination: Recovering reflectance models
 of real scenes from photographs. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, 1999. 7
 702
- [79] Ye Yu, Abhimitra Meka, Mohamed Elgharib, Hans-Peter Seidel, Christian Theobalt, and William A. P. Smith. Selfsupervised outdoor scene relighting. In *ECCV*, 2020. 6
- [80] Jinsong Zhang, Kalyan Sunkavalli, Yannick Hold-Geoffroy, Sunil Hadap, Jonathan Eisenman, and Jean-François Lalonde. All-weather deep outdoor lighting estimation. In *CVPR*, 2019.
- [81] Kai Zhang, Fujun Luan, Qianqian Wang, Kavita Bala, and Noah Snavely. Physg: Inverse rendering with spherical gaussians for physics-based material editing and relighting. In *CVPR*, 2021. 8
- [82] Kai Zhang, Fujun Luan, Zhengqi Li, and Noah Snavely. Iron: Inverse rendering by optimizing neural sdfs and materials from photometric images. In *CVPR*, 2022. 7
- [83] Ruo Zhang, Ping-Sing Tsai, James Edwin Cryer, and Mubarak Shah. Shape from shading: A survey. *IEEE TPAMI*, 1999. 7
 718
- [84] Shuyang Zhang, Runze Liang, and Miao Wang. Shadowgan: Shadow synthesis for virtual objects with conditional adversarial networks. *Computational Visual Media*, 2019. 9
- [85] Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul Debevec, William T Freeman, and Jonathan T Barron. Nerfactor: Neural factorization of shape and reflectance under an unknown illumination. ACM TOG, 2021. 2, 4, 8
- [86] Yuxuan Zhang, Wenzheng Chen, Huan Ling, Jun Gao, Yinan Zhang, Antonio Torralba, and Sanja Fidler. Image gans meet differentiable rendering for inverse graphics and interpretable 3d neural rendering. arXiv, 2020. 7
- [87] Yuanqing Zhang, Jiaming Sun, Xingyi He, Huan Fu, Rongfei Jia, and Xiaowei Zhou. Modeling indirect illumination for inverse rendering. In *CVPR*, 2022. 2, 4, 8, 9
- [88] Tinghui Zhou, Philipp Krahenbuhl, and Alexei A Efros. Learning data-driven reflectance priors for intrinsic image decomposition. In *ICCV*, 2015. 7