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

Anonymous CVPR submission

Paper ID



Input video

Inverse rendering

Figure 1. We present UrbanIR (Urban Scene Inverse Rendering), a novel, realistic, and relightable neural scene model. UrbanIR concurrently infers shape, albedo, visibility, and more from a single video of large-scale, unbounded scenes. The resulting representation facilitates realistic and controllable editing, delivering photorealistic free-viewpoint renderings (last column) of relit scenes (top row), inserted objects (middle row), and nighttime simulation (bottom row).

Abstract

001 We present UrbanIR (Urban Scene Inverse Rendering), a new inverse graphics model that enables realistic, free-002 viewpoint renderings of scenes under various lighting condi-003 004 tions with a single video. It accurately infers shape, albedo, visibility, and sun and sky illumination from wide-baseline 005 videos, such as those from car-mounted cameras, differ-006 007 ing from NeRF's dense view settings. In this context, standard methods often yield subpar geometry and material 800 009 estimates, such as inaccurate roof representations and nu-010 merous 'floaters'. UrbanIR addresses these issues with novel 011 losses that reduce errors in inverse graphics inference and 012 rendering artifacts. Its techniques allow for precise shadow 013 volume estimation in the original scene. The model's outputs 014 support controllable editing, enabling photorealistic free-015 viewpoint renderings of night simulations, relit scenes, and 016 inserted objects, marking a significant improvement over existing state-of-the-art methods. Our code and data will be 017 018 made publicly available upon acceptance.

1. Introduction 019

We show how to build a model that allows realistic, free-020 021 viewpoint renderings of a scene under novel lighting conditions from a video. So, for example, a sunny afternoon video 022 of a large urban scene can be shown at different times of day 023 or night (as in Fig. 1), viewed from novel viewpoints, and 024 shown with inserted objects. Our method - UrbanIR (Ur-025 ban Scene Inverse Rendering) - computes an inverse graph-026 ics representation from the video. UrbanIR jointly infers 027 shape, albedo, visibility, and sun and sky illumination from 028 a single video of unbounded outdoor scenes with unknown 029 lighting. The resulting representations enable controllable 030 editing, delivering photorealistic free-viewpoint renderings 031 of relit scenes and inserted objects, as demonstrated in Fig. 1. 032

CVPR

#

UrbanIR obtains its intrinsic scene representations from 033 a video under a single illumination condition, but producing 034 realistic novel views requires accurate inferences of physi-035 cal parameters. UrbanIR uses a novel visibility rendering 036 scheme and loss to make precise estimates of shadow vol-037 umes in the original scene and so control albedo errors. Ur-038 banIR combines monocular intrinsic decomposition and in-039 verse rendering with other key contributions to control errors 040 in renderings. To our knowledge, UrbanIR is the first in its 041 class capable of performing inverse rendering and relighting 042 applications from a single monocular video in large-scale 043 scenes, without requiring multiple illumination conditions, 044 depth sensing, or both. 045

UrbanIR representations are constructed from cameras 046

090

091

092

093

094

095

096

097

102

114



Figure 2. **Rendering Pipeline.** UrbanIR retrieves scene intrinsics (normal N, semantics S, albedo A) from camera rays, and estimate visibility V from tracing rays to the light source. The shading model computes diffuse and specular reflection and adds ambient sky light \mathbf{L}_{sky} for the final shading map. We multiply shading & albedo, and render the sky appearance for final rendering.

047 mounted on cars with a narrow range of views of each scene point. Typical NeRF-style systems yield poor geometry es-048 049 timates (for example, roofs) and numerous "floaters" under these conditions; they are usually trained with a wide range 050 051 of views. Our experiments showcase that UrbanIR outper-052 forms these baselines with significantly reduced artifacts 053 in our sparse view setting. Finally, we also show how to 054 use UrbanIR to simulate night scenes from a single daytimecaptured video, producing a controllable, realistic, physically 055 056 plausible, and consistent simulation. In summary, our key contributions are: 057

- We present UrbanIR for recovering a *relightable* neural radiance field in a constrained setting of an *unbounded scene*, using a *single monocular video* captured under a *single illumination condition*.
- We describe a novel inverse rendering framework that *builds precise shadow volumes* in large outdoor scenes with heavy shadows, resulting in significant improvements in inverse graphics estimates and relighting.
- We demonstrate a new physics-informed night simulation framework. To our knowledge, UrbanIR is the first simulation to offer realistic, *free-viewpoint night simulation* from a single daytime video capture.

070 2. Method

The illustration of the rendering pipeline is shown in Fig. 2.
Please refer to our supplementary materials for the complete
method description.

3. Experiment Results

075 **3.1.** Datasets

We evaluate UrbanIR on two datasets: the KITTI-360
dataset [3] and the Waymo Open Dataset [8]. The KITTI360 dataset [3] consists of 9 stereo video sequences showcasing urban scenes. For our analysis, we selected 7 nonoverlapping clipped sequences, each containing around 100
images. These sequences cover various light directions, vehi-

cle trajectories, and layouts of buildings and vegetation. The 082 dataset includes RGB images from stereo cameras, semantic 083 labels, camera poses, and RTK-GPS poses. On the other 084 hand, the Waymo Open Dataset (WOD) [8] captures driving 085 sequences from five cameras and one 64-beam LiDAR sen-086 sor at 10 Hz. However, we only used the single camera from 087 the front view and did not use any LiDAR information for 088 our evaluation. 089

Quantitative evaluation of relighting sequences is difficult as most datasets only capture the same location under a single illumination, and no ground truth for relighting is available. Therefore, we recorded a scene at different times of the day, covering different illuminations. The images were captured by a stereo camera, and the poses were estimated using RTK-GPS information.

3.2. Baselines

We compare UrbanIR with scene relighting and editing methods:098ods: FEGR [9], Colmap MVS [7], Instruct NeRF2NeRF [2],099NeRF-OSR [5], RelightNet [10]. Implementation details are100in supplementary.101

3.3. Decomposition Quality

We evaluate intrinsic decomposition on the Waymo Open 103 Dataset [8] and present the comparison in Fig 3. NeRF-104 OSR [4] requires multi-illumination as input and fails to 105 decompose albedo and shadow, leaving severe artifacts due 106 to noisy normal estimation. FEGR [9] uses five cameras and 107 LiDAR for reconstruction but still bakes shadow patterns 108 into the albedo and normal. However, UrbanIR only requires 109 a single camera as input without any LiDAR information. By 110 integrating monocular prior in the optimization process, it 111 successfully decomposes clean albedo, normal, and shadow 112 maps under single illumination. 113

3.4. Relighting Quality

Relighting under various sunlight conditions is evaluated in 115 Fig. 5, ??. NeRF-OSR [5] cannot simulate shadows under 116 novel light conditions. While Blender [1] can change the 117 lighting parameters explicitly, they either cast bad shadows 118 due to incomplete geometry or do not cast new shadows 119 at all; further, the original shadow remains unchanged in 120 the image. We implement a mesh-based visibility baseline 121 which extracts mesh with marching cubes for visibility cal-122 culation. It generates different shadows according to light 123 conditions, but the mesh on the edges of and outside the 124 training views is poor because there are few observations. 125 This leads to noisy geometry and incomplete shadows on 126 the ground. In contrast, UrbanIR synthesizes sharp shadows 127 and varying surface shading following the sun's direction. 128 Further, the original scene shadows are largely absent. This 129 allows synthesizing images at night (Fig. ??) by inserting 130 car headlights and streetlights, without distracting effects 131

CVPR 2024 Submission #. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.



Figure 3. Intrinsic Decomposition of Waymo Open Dataset [8]. FEGR does not release code, so we directly use the images from their paper, and the shadow map of this viewpoint is not available.



Figure 4. Controllable Relighting of Waymo Open Dataset [8]. The first row shows different lighting during a day, and the second row changes the input image into night-time with different lighting configurations.

from the original shadows. Moreover, the relighting results
obtained from UrbanIR are highly controllable, as demonstrated in Fig. 4. Different light directions and intensities
were used to adjust the relighting outcomes. Light sources
were also added and turned on and off.

137 3.5. Quantitative Evaluation

138 The results of the quantitative evaluation can be found in Tab. 1. We tested the novel view synthesis on KITTI-360 [3] 139 using 10 images as the novel views for all 7 sequences. Ur-140 banIR is better than NeRF-OSR in all metrics, indicating that 141 142 our model can decompose intrinsic well and produce highquality images. To evaluate the relighting in novel views, 143 we captured videos of the same scene in the morning and 144 afternoon. After optimizing models at both sequences in-145 dividually, we performed relighting by exchanging lighting 146 147 parameters and calculating image metrics with the ground 148 truth capture. Our method also outperformed NeRF-OSR in

all metrics significantly. The qualitative results can be seen149in Fig. 7. UrbanIR was successful in removing existing shadows, changing the shading on the building, and modifying150the sky texture during different times of the day.151

3.6. Object Insertion

153

162

The object insertion pipeline is described in Sect. ??. In 154 Fig. 6, we compare with baselines by inserting a yellow 155 cube and moving along the road. Mesh+Blender cannot 156 synthesize complete geometry and shadow. Without visibil-157 ity optimization in row (B), the scene shadow on objects is 158 noisy. Our complete model has an accurate shadow volume 159 so that shadows cast on the object by the environment are 160 well represented. 161

4. Discussion

In this work, we investigated the task of inverse rendering of unbounded outdoor scenes under single illumination. This



Figure 5. **Rendering and relighting comparison**. We show a set of two scenes comparing different methods. Each column shows a different sun position, the first column showing original images. For each set and from top to bottom, we have (a) COLMAP [7] + Blender [1] (b) Mesh-based visibility (NeRF-Mesh) and (c) UrbanIR (Ours). In COLMAP, Shadows are present in the scene; however, they are "baked-in" and cannot be manipulated or relit independently. For NeRF-Mesh, weak observations of the scene geometry result in poor quality visibility estimation during direct observation-based reconstruction. In contrast, for our approach, UrbanIR, our visibility optimization enables realistic and controllable relighting effects.



Figure 6. **Dynamic Object Insertion with Shadow Volume.** Our method produces accurate estimates of shadow volumes where others cannot. This can be visualized by inserting a simple object into the scene, and then looking at shadows cast onto that object. (A) COLMAP dense reconstruction [6, 7] + Blender [1] (B) Ours without visibility optimization and (C) Ours with visibility optimization improves shadows.



B Ground Truth (3pm) B Model + B Light (Ours) A Model + B Light (Ours) A Model + B Light (NeRF-OSR) Figure 7. Novel view and novel light synthesis.

task is ill-posed and extremely challenging due to the sparsity of observations across space and time. To overcome
this challenge and successfully decompose various scene
intrinsic properties, we utilized prior knowledge such as pretrained networks and regularization to reduce the uncertainty

space and improve the performance of downstream applica-170 tions like relighting and object insertion. However, there are 171 limitations. Our optimization process can be affected by the 172 noisy predictions from prior models and requires careful tun-173 ing of our losses. Sometimes, shadows cannot be removed 174 entirely in the albedo field, and they may still appear in the 175 final images. Additionally, the visibility optimization refines 176 only the geometry along the light direction, which means 177 that large changes in the sun's direction can lead to poor 178 shadows when the geometry estimates are not accurate. 179

References

- Blender Online Community. *Blender a 3D modelling and rendering package*. Blender Foundation, Stichting Blender
 Foundation, Amsterdam, 2018. 2, 4
- [2] Ayaan Haque, Matthew Tancik, Alexei Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-nerf2nerf: Editing 3d scenes with instructions. In *Proceedings of the IEEE/CVF*186

180

187

195

196

197 198

199

200 201

202 203 International Conference on Computer Vision, 2023. 2

- [3] 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. 2, 3
- 191 [4] Viktor Rudnev, Mohamed Elgharib, William Smith, Lingjie
 192 Liu, Vladislav Golyanik, and Christian Theobalt. Neural
 193 radiance fields for outdoor scene relighting. *ECCV*, 2022. 2,
 194 3, 4
 - [5] 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. 2, 4
 - [6] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 4
 - [7] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *ECCV*, 2016. 2, 4
- [8] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien 204 205 Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, 206 Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, 207 Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott Et-208 tinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang, 209 Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. 210 Scalability in perception for autonomous driving: Waymo 211 open dataset. In CVPR, 2020. 2, 3
- [9] 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.
 216
 2, 3
- [10] Ye Yu, Abhimitra Meka, Mohamed Elgharib, Hans-Peter
 Seidel, Christian Theobalt, and William A. P. Smith. Selfsupervised outdoor scene relighting. In *ECCV*, 2020. 2