

NeRV: Neural Reflectance and Visibility Fields for Relighting and View Synthesis

Pratul P. Srinivasan
Google Research

Boyang Deng
Google Research

Xiuming Zhang
MIT

Matthew Tancik
UC Berkeley

Ben Mildenhall
UC Berkeley

Jonathan T. Barron
Google Research

Abstract

We present a method that takes as input a set of images of a scene illuminated by unconstrained known lighting, and produces as output a 3D representation that can be rendered from novel viewpoints under arbitrary lighting conditions. Our method represents the scene as a continuous volumetric function parameterized as MLPs whose inputs are a 3D location and whose outputs are the following scene properties at that input location: volume density, surface normal, material parameters, distance to the first surface intersection in any direction, and visibility of the external environment in any direction. Together, these allow us to render novel views of the object under arbitrary lighting, including indirect illumination effects. The predicted visibility and surface intersection fields are critical to our model’s ability to simulate direct and indirect illumination during training, because the brute-force techniques used by prior work are intractable for lighting conditions outside of controlled setups with a single light. Our method outperforms alternative approaches for recovering relightable 3D scene representations, and performs well in complex lighting settings that have posed a significant challenge to prior work.

1. Introduction

A central problem in computer vision is that of inferring the physical geometry and material properties that together explain observed images. In addition to its importance for recognition and robotics, a solution to this open problem would have significant value for computer graphics — the ability to create realistic 3D models from standard photos could democratize 3D content creation and allow anyone to use real-world objects in photography, filmmaking, and game development. In this paper, we work towards this goal and present an approach for estimating a volumetric 3D representation from images of a scene under arbitrary known lighting conditions, such that high-quality novel images can be rendered from arbitrary unseen viewpoints and under

Figure 1: We optimize a Neural Reflectance and Visibility Field (NeRV) 3D representation from a set of input images of a scene illuminated by known but unconstrained lighting. Our NeRV representation can be rendered from novel views under arbitrary lighting conditions not seen during training. Here, we visualize example input data and renderings for two scenes. The first two output rendered images for each scene are from the same viewpoint, each illuminated by a point light at a different location, and the last image is from a different viewpoint under a random colored illumination.

novel unobserved lighting conditions, as shown in Figure 1.

The vision and graphics research communities have recently made substantial progress towards the novel view synthesis portion of this goal. The Neural Radiance Fields (NeRF) [32] approach has shown that it is possible to synthesize photorealistic images of scenes by training a simple neural network to map 3D locations in the scene to a continuous field of volume density and color. Volume rendering is trivially differentiable, so the parameters of a NeRF can be optimized for a single scene by using gradient descent to minimize the difference between renderings of the NeRF and a set of observed images. Though NeRF produces compelling results for view synthesis, it does not provide a solution for relighting. This is because NeRF models just the amount of

outgoing light from a location — the fact that this outgoing light is the result of interactions between incoming light and the material properties of an underlying surface is ignored.

At first glance, extending NeRF to enable relighting appears to require only changing the image formation model: instead of modeling scenes as fields of density and view-dependent color, we can model surface normals and material properties (*e.g.* the parameters of a bi-directional reflectance distribution function (BRDF)) and simulate the transport of the scene’s light sources (which we assume are known) according to the rules of physically based rendering [38]. However, simulating the attenuation and reflection of light by particles is fundamentally challenging in NeRF’s neural volumetric representation because content can exist *anywhere* within the scene, and determining the density at any location requires querying a neural network. Consider the naïve procedure for computing the radiance along a single camera ray due to direct illumination, as illustrated in Figure 2: First, we query NeRF’s multi-layer perceptron (MLP) for the volume density at samples along the camera ray to determine the amount of light reflected by particles at each location that reaches the camera. For each location along the camera ray, we then query the MLP for the volume density at densely-sampled points between the location and every light source to estimate the attenuation of light before it reaches that location. This procedure quickly becomes prohibitively expensive if we want to model environment light sources or global illumination, in which case scene points may be illuminated from all directions. Prior methods for estimating relightable volumetric representations from images have not overcome this challenge, and can only simulate direct illumination from a single point light source when training.

The problem of efficiently computing visibility is well explored in the graphics literature. In standard raytracing graphics pipelines, where the scene geometry is fixed and known ahead of time, a common solution is to precompute a data structure that can be efficiently queried to obtain the visibility between pairs of scene points, or between scene points and light sources. This can be accomplished with approaches including octrees [44], distance transforms [8], or bounding volume hierarchies [38]. But these existing approaches do not provide a solution to our task — our geometry is unknown, and our model’s estimate of geometry changes constantly as it is optimized. Though conventional data structures could perhaps be used to accelerate rendering *after* optimization is complete, we need to efficiently query the visibility between points *during* optimization, and existing solutions are prohibitively expensive to rebuild after each training iteration (of which there may be millions).

In this work, we present a method to train a NeRF-like model that can simulate realistic environment lighting and global illumination. Our key insight is to train an MLP to act as a lookup table into a *visibility field* during rendering.

Figure 2: We visualize how “Neural Visibility Fields” reduce the computational burden of volume rendering a camera ray with direct (top) and one-bounce indirect (bottom) illumination compared to naïve raymarching, alongside each solution’s computational complexity (n is the number of samples along each ray, l is the number of light sources, and d is the number of sampled indirect illumination directions). Black dots represent evaluating a shape MLP for volume density at a position, red arrows represent evaluating the visibility MLP at a position along a direction, and the blue arrow represents evaluating the visibility MLP for the expected termination depth of a ray at a position along a direction (output visibility multipliers and termination depths from the visibility MLP are displayed as text). Brute-force light transport simulation through NeRF’s volumetric representation with naïve raymarching (left) is intractable. By approximating visibility with a neural visibility field (right) that is optimized alongside the shape MLP, we are able to make optimization with complex illumination tractable.

Instead of estimating light or surface visibility at a given 3D position along a given direction by densely evaluating an MLP for the volume density along the corresponding ray (which would be prohibitively expensive), we simply query this visibility MLP to estimate visibility and expected termination depth in any direction (see Figure 2). This visibility MLP is optimized alongside the MLP that represents volume density, and is supervised to be consistent with the volume density samples observed during optimization. Using this neural approximation of the true visibility field significantly eases the computational burden of estimating volume rendering integrals while training. Our resulting system, which we call “NeRV” (“Neural Reflectance and Visibility Fields”) enables the recovery of a NeRF-like model that supports relighting in addition to view synthesis. While previous solutions for relightable NeRFs [3] were limited to controlled settings which required input images to be illuminated by a single point light, NeRV supports training with arbitrary environment lighting and “one-bounce” indirect illumination.

2. Related Work

Neural Radiance Fields [32] can be thought of as a modern neural reformulation of the classic problem of scene reconstruction: given multiple images of a scene, inferring the underlying geometry and appearance that best explains those images. While classic approaches have largely relied on discrete representations such as textured meshes [16, 53] and voxel grids [48], NeRF has demonstrated that a continuous volumetric function, parameterized as an MLP, is able to represent complex scenes and render photorealistic novel views. NeRF works well for view synthesis, but it does not enable relighting because it has no mechanism to disentangle the outgoing radiance of a surface into an incoming radiance and an underlying surface material.

This problem of attributing what aspects of an image are due to material, lighting, or geometric variation is commonly referred to as “intrinsic image estimation” [2, 22] or “inverse rendering” [40, 46], and is a classic problem in computer vision and graphics. These classical approaches have derived insightful observations about separating a single image into shading and reflectance components [18], inferring surface normals from the appearance of an image’s shading [17], or jointly inferring shape, illumination, and reflectance from a single image [1], but they are not designed to recover full 3D models that can be used for graphics applications.

The difficulty of this problem (a consequence of its underconstrained nature) is typically addressed using one of the following strategies: 1) learning priors on shape, illumination, and reflectance, 2) assuming known geometry, or 3) using multiple input images of the scene under different lighting conditions. Most recent single-image inverse rendering methods [23, 25, 45, 49, 56, 60] belong to the first category, and use large datasets of images with labeled geometry and materials to train convolutional neural networks to predict these properties. Most prior works in inverse rendering that recover full 3D models for graphics applications [57] fall under the second category, and use 3D geometry obtained from active scanning [37, 47, 61], proxy models [7, 10, 11], silhouette masks [35, 58], or multiview stereo [33] as a starting point before recovering reflectance and refined geometry.

Our method belongs to the third category; we only require posed input images of a scene under different (known) lighting conditions. The most related prior works are Deep Reflectance Volumes [4], which estimates voxel geometry and BRDF parameters, and the follow-up work Neural Reflectance Fields [3], which replaces Deep Reflectance Volume’s voxel grid with a continuous volume represented by an MLP. Our work extends Neural Reflectance Fields (which requires scenes to only be illuminated by a single point light at a time due to their brute-force visibility computation strategy visualized in Figure 2, and only models direct illumination) to work for arbitrary lighting and global illumination.

We also take inspiration from recent works that replace

discrete voxel and mesh geometry representations with MLPs that approximate a continuous 3D function by mapping from an input 3D location to scene properties at that location. This strategy has been explored for the tasks of representing shapes [9, 12, 30, 36, 50, 54] and scenes under fixed lighting for view synthesis [26, 32, 34, 51, 59]. One technique that has been used for relighting these neural representations is to condition the MLP’s output appearance on a latent code that encodes a per-image lighting, as in NeRF in the Wild [28] (as well as previously with discretized scene representations [24, 31]). Although this strategy can effectively explain the appearance variation of training images, it cannot be used to render the same scene under new lighting conditions not observed during training (Figure 6) because it does not utilize the physics of light transport.

Our method is inspired by a long line of work in graphics that explores precomputation [42, 52] and approximation [6, 14, 41, 43] strategies to efficiently compute global illumination in physically-based rendering. Our “Neural Visibility Fields” can be thought of as a neural analogue to visibility precomputation techniques, and is specifically designed for use in our neural inverse rendering setting where geometry is dynamically changing during optimization.

3. Method

We extend NeRF to include the simulation of light transport, which allows NeRFs to be rendered under arbitrary novel illumination conditions. Instead of modeling a scene as a continuous 3D field of particles that absorb and *emit* light as in NeRF, we represent a scene as a 3D field of *oriented* particles that absorb and *reflect* the light emitted by external light sources (Section 3.2). Naïvely simulating light transport through this model is inefficient and unable to scale to simulate realistic lighting conditions or global illumination. We remedy this by introducing a neural visibility field representation (optimized alongside NeRF’s volumetric representation) that allows us to efficiently query the point-to-light and point-to-point visibilities needed to simulate light transport (Section 3.3). The resulting Neural Reflectance and Visibility Field (NeRV) is visualized in Figure 3.

3.1. NeRF Overview

NeRF represents a scene as a continuous function, parameterized by a “radiance” MLP whose input is a 3D position and viewing direction, and whose output is the volume density and radiance L_e (RGB color) emitted by particles at that location along that viewing direction. NeRF uses standard emission-absorption volume rendering [19] to compute the observed radiance $L(c, \omega_o)$ (the rendered pixel color) at camera location c along direction ω_o as the integral of the product of three quantities at any point $x(t) = c - t \omega_o$ along the ray: the visibility $V(x(t), c)$, which indicates the fraction of emitted light from position $x(t)$ that reaches the

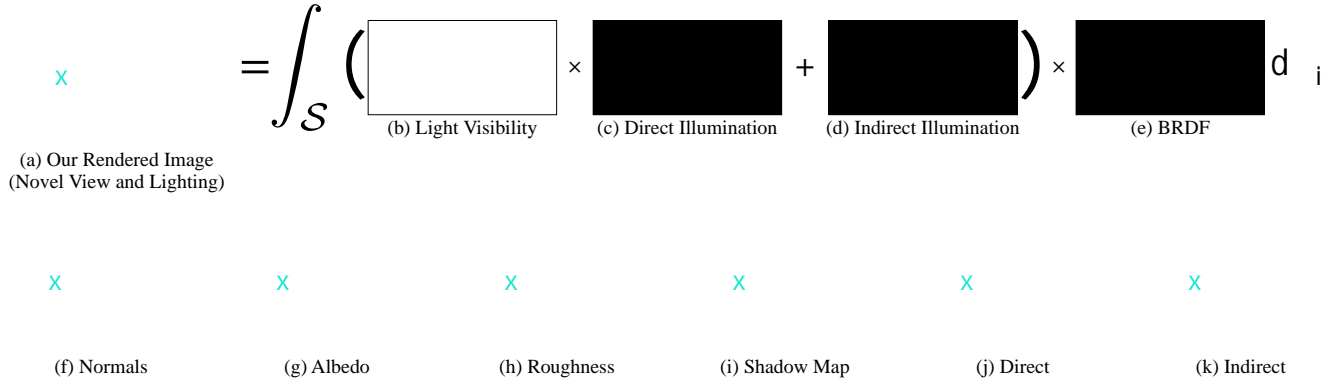


Figure 3: (a) Given any continuous 3D location as input, such as the point at the cyan “x”, NeRV outputs the volume density as well as: (b) The visibility to a spherical environment map surrounding the scene, which is multiplied by (c) the direct illumination at that point and added to (d) the estimated indirect illumination map at that point to determine the full incident illumination. This is then multiplied by (e) the predicted BRDF, and integrated over all incoming directions to determine the outgoing radiance at that point. In the bottom row, we visualize these outputs for the full rendered image: (f) Surface normals, and BRDF parameters for (g) diffuse albedo and (h) specular roughness. We can use the predicted visibilities to compute the fraction of the total illumination that is actually incident at any location, visualized as (i) a shadow map. We also show the same rendered viewpoint if it were lit by only (j) direct and (k) indirect illumination.

camera at c , the density $\rho(x(t))$, and the emitted radiance $L_e(x(t), \omega_o)$ along the viewing direction ω_o :

$$L(c, \omega_o) = \int_0^t V(x(t), c) \rho(x(t)) L_e(x(t), \omega_o) dt, \quad (1)$$

$$V(x(t), c) = \exp \left(- \int_0^t \rho(x(s)) ds \right). \quad (2)$$

A NeRF is recovered from observed input images of a scene by sampling a batch of observed pixels, sampling the corresponding camera rays of those pixels at stratified random points to approximate the above integral using numerical quadrature [29], and optimizing the weights of the radiance MLP via gradient descent to minimize the error between the estimated and observed pixel colors.

3.2. Neural Reflectance Fields

A NeRF representation does not separate the effect of incident light from the material properties of surfaces. This means that NeRF is only able to render views of a scene under the fixed lighting conditions presented in the input images — a NeRF cannot be relit. Modifying NeRF to enable relighting is straightforward, as initially demonstrated by the Neural Reflectance Fields work of Bi *et al.* [3]. Instead of representing a scene as a field of particles that emit light, it is represented as a field of particles that reflect incoming light. With this, given an arbitrary lighting condition, we can simulate the transport of light through the volume as it is reflected by particles until it reaches the camera with a

standard volume rendering integral [19]:

$$L(c, \omega_o) = \int_0^t V(x(t), c) \rho(x(t)) L_r(x(t), \omega_o) dt, \quad (3)$$

$$L_r(x, \omega_o) = \int_S L_i(x, \omega_i) R(x, \omega_i, \omega_o) d\omega_i, \quad (4)$$

where the view-dependent emission term $L_e(x, \omega_o)$ in Equation 1 is replaced with an integral over the sphere S of incoming directions, of the product of the incoming radiance L_i from any direction and a reflectance function R (often called a phase function in volume rendering) which describes how much light arriving from direction ω_i is reflected towards direction ω_o . We follow Bi *et al.* and use the standard microfacet BRDF described by Walter *et al.* [55] as the reflectance function, so R at any 3D location is parameterized by a diffuse RGB albedo, a scalar specular roughness, and a surface normal. We replace NeRF’s radiance MLP with two MLPs: a “shape” MLP that outputs volume density and a “reflectance” MLP that outputs BRDF parameters (3D diffuse albedo a and 1D roughness r) for any input 3D point: $\text{MLP}_\rho : x \rightarrow \rho$, $\text{MLP}_R : x \rightarrow (a, r)$. Instead of parameterizing the 3D surface normal n as a normalized output of the shape MLP, as in Bi *et al.* [3], we compute n as the negative normalized gradient vector of the shape MLP’s output with respect to x , computed using automatic differentiation. We further discuss this choice in Section 4.2.

3.3. Light Transport via Neural Visibility Fields

Although modifying NeRF to enable relighting is straightforward, estimating the volume rendering integral for general lighting scenarios is computationally challenging with a continuous volumetric representation such as NeRF. Fig-

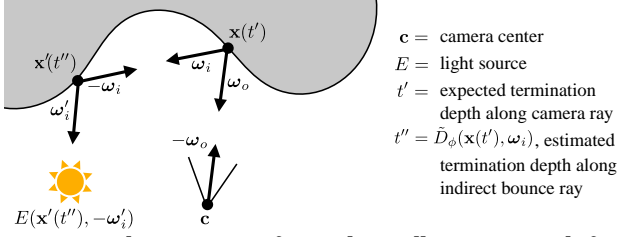


Figure 4: The geometry of an indirect illumination path from camera to light source, and a visualization of our notation.

Figure 2 visualizes the scaling properties that make simulating volumetric light transport particularly difficult. Even if we only consider direct illumination from light sources to a scene point, a brute-force solution is already challenging for more than a single point light source as it requires repeatedly querying the shape MLP for volume density along paths from each scene point to each light source. Moreover, general scenes can be illuminated by light arriving from *all* directions, and addressing this is imperative to recovering relightable representations in unconstrained scenarios. Simulating even simple global illumination in a brute-force manner is intractable: rendering a *single ray* in our scenes under one-bounce indirect illumination with brute-force sampling would require a *petaflop* of computation, and we need to render roughly a *billion* rays over the course of training.

We ameliorate this issue by replacing several brute-force volume density integrals with learned approximations. We introduce a “visibility” MLP that emits an approximation of the environment lighting visibility at any input location along any input direction, as well as an approximation of the expected termination depth of the corresponding ray: MLP : $(x, \omega) \rightarrow (\tilde{V}, \tilde{D})$. When rendering, we use these MLP-approximated quantities in place of their actual values:

$$V(x, \omega) = \exp \left(- \int_0^\infty \rho(x + s\omega) ds \right), \quad (5)$$

$$D(x, \omega) = \int_0^\infty \exp \left(- \int_0^t \rho(x + s\omega) ds \right) \rho(x + t\omega) dt. \quad (6)$$

In Section 3.5 we place losses on the visibility MLP outputs (\tilde{V}, \tilde{D}) to encourage them to resemble the (V, D) corresponding to the current state of the shape MLP.

Below, we provide a detailed walkthrough of how our Neural Visibility Field approximations simplify the volume rendering integral computation. Figure 4 is provided for reference. We first decompose the reflected radiance $L_r(x, \omega_o)$ into its direct and indirect illumination components. Let us define $L_e(x, \omega_i)$ as radiance due to a light source arriving at point x from direction ω_i . As defined in Equation 3, $L(x, \omega_i)$ is the estimated incoming radiance at location x from direction ω_i . This means the incident illumination L_i decomposes into $L_e + L$ (direct plus indirect lighting). The

shading calculation for L_r then becomes:

$$\begin{aligned} L_r(x, \omega_o) &= \int_S (L_e(x, \omega_i) + L(x, -\omega_i)) R(x, \omega_i, \omega_o) d\omega_i \quad (7) \\ &= \underbrace{\int_S L_e(x, \omega_i) R(x, \omega_i, \omega_o) d\omega_i}_{\text{component due to direct lighting}} + \underbrace{\int_S L(x, -\omega_i) R(x, \omega_i, \omega_o) d\omega_i}_{\text{component due to indirect lighting}}. \end{aligned}$$

To calculate incident direct lighting L_e we must account for the attenuation of the (known) environment map E due to the volume density along the incident illumination ray ω_i :

$$L_e(x, \omega_i) = V(x, \omega_i) E(x, -\omega_i). \quad (8)$$

Instead of evaluating V as another line integral through the volume, we use the visibility MLP’s approximation \tilde{V} . With this, our full calculation for the direct lighting component of camera ray radiance $L(c, \omega_o)$ simplifies to:

$$\int_0^\infty V(x(t), c - x(t)) \int_S \tilde{V}(x(t), \omega_i) E(x(t), -\omega_i) R(x(t), \omega_i, \omega_o) d\omega_i dt. \quad (9)$$

By approximating the integrals along rays from each point on the camera ray toward each environment direction when computing the color of a pixel due to direct lighting, we have reduced the complexity of rendering with direct lighting from quadratic in the number of samples per ray to linear. Next, we focus on the more difficult task of accelerating the computation of rendering with indirect lighting, for which a brute force approach would scale cubically with the number of samples per ray. We make two approximations to reduce this intractable computation. Our first approximation is to replace the outermost integral (the accumulated radiance reflected towards the camera at each point along the ray) with a single point evaluation by treating the volume as a hard surface located at the expected termination depth $t = D(c, -\omega_o)$. Note that we do not use the visibility MLP’s approximation of t here, since we are already sampling (x) along the camera ray. This reduces the indirect contribution of $L(c, \omega_o)$ to a spherical integral at a single point $x(t)$:

$$L(x(t), -\omega_i) R(x(t), \omega_i, \omega_o) d\omega_i. \quad (10)$$

To simplify the recursive evaluation of L inside this integral, we limit the indirect contribution to a single bounce, and use the hard surface approximation a second time to replace the integral along a ray for each incoming direction:

$$L(x(t), -\omega_i) = \int_S L_e(x(t), \omega_j) R(x(t), \omega_j, -\omega_i) d\omega_j, \quad (11)$$

where $t = \tilde{D}(x(t), \omega_i)$ is the expected intersection depth along the ray $x(t) = x(t) + t\omega_i$ as approximated by the visibility MLP. Thus the expression for the component of camera ray radiance $L(c, \omega_o)$ due to indirect lighting is:

$$\int_S L_e(x(t), \omega_j) R(x(t), \omega_j, -\omega_i) d\omega_j R(x(t), \omega_i, \omega_o) d\omega_i, \quad (12)$$

and fully expanding the direct radiance $L_e(x(t), \omega_i)$ incident at each secondary intersection point gives us:

$$\tilde{V}(x(t), \omega_i) E(x(t), -\omega_i) R(x(t), \omega_i, \omega_o) d\omega_i \quad (13)$$

Finally, we can write out the complete volume rendering equation used by NeRV as the sum of Equations 9 and 13:

$$L(c, \omega_o) = \int_0^{\infty} V(x(t), c - x(t)) \tilde{V}(x(t), \omega_i) E(x(t), -\omega_i) R(x(t), \omega_i, \omega_o) dt + \int_S \tilde{V}(x(t), \omega_i) E(x(t), -\omega_i) R(x(t), \omega_i, \omega_o) d\omega_i \quad (14)$$

Figure 2 illustrates how the approximations made by NeRV reduce the computational complexity of computing direct and indirect illumination from quadratic and cubic (respectively) to linear. This enables the simulation of direct illumination from environment lighting and one-bounce indirect illumination within the training loop of optimizing a continuous relightable volumetric scene representation.

3.4. Rendering

To render a camera ray $x(t) = c - t \omega_o$ passing through a NeRV, we estimate the volume rendering integral in Equation 14 using the following procedure:

- 1) We draw 256 stratified samples along the ray and query the shape and reflectance MLPs for the volume densities, surface normals, and BRDF parameters at each point: $\rho = \text{MLP}(x(t))$, $n = x(t) \times \text{MLP}(x(t))$, $(a, b) = \text{MLP}(x(t))$.
- 2) We shade each point along the ray with direct illumination by estimating the integral in Equation 9. First, we generate $E(x(t), -\omega_i)$ by sampling the known environment lighting on a 12×24 grid of directions ω_i on the sphere around each point. We then multiply this by the predicted visibility $\tilde{V}(x(t), \omega_i)$ and microfacet BRDF values $R(x(t), \omega_i, \omega_o)$ at each sampled ω_i , and integrate this product over the sphere to produce the direct illumination contribution.
- 3) We shade each point along the ray with indirect illumination by estimating the integral in Equation 13. First, we compute the expected camera ray termination depth $t = D(c, -\omega_o)$ using the density samples from Step 1. Next, we sample 128 random directions on the upper hemisphere from $x(t)$ and query the visibility MLP for the expected termination depths along each of these rays $t = D(x(t), \omega_i)$ to compute the secondary surface intersection points $x(t) = x(t) + t \omega_i$. We then shade each of these points with direct illumination by following the procedure in Step 2. This estimates the indirect illumination incident at $x(t)$, which we then multiply by the microfacet BRDF values $R(x(t), \omega_i, \omega_o)$ and integrate over the sphere to produce the indirect illumination contribution.
- 4) The total reflected radiance at each point along the camera ray $L_r(x(t), \omega_o)$ is the sum of the quantities from Steps 2

and 3. We composite these along the ray to compute the pixel color using the same quadrature rule [29] used in NeRF:

$$L(c, \omega_o) = \int_0^{\infty} V(x(t), c - x(t)) L_r(x(t), \omega_o) dt \quad (15)$$

$$V(x(t), c) = \exp - \int_{s < t} \rho(x(s)) ds, \quad (z) = 1 - \exp(-z),$$

where Δ is the distance between samples along the ray.

3.5. Training and Implementation Details

Instead of directly passing 3D coordinates x and direction vectors ω to the MLPs, we map these inputs using NeRF’s positional encoding [32, 54], with a maximum frequency of 2^7 for 3D coordinates and 2^4 for 3D direction vectors. The shape and reflectance MLPs each use 8 fully-connected ReLU layers with 256 channels. The visibility MLP uses 8 fully-connected ReLU layers with 256 channels each to map the encoded 3D coordinates x to an 8-dimensional feature vector which is concatenated with the encoded 3D direction vector ω and processed by 4 fully-connected ReLU layers with 128 channels each.

We train a separate NeRV representation from scratch for each scene, which requires a set of posed RGB images and corresponding lighting environments. At each training iteration we randomly sample a batch of 512 pixel rays R from the input images and use the previously-described procedure to render these pixels from the current NeRV model. We additionally sample 256 random rays R per training iteration that intersect the volume, and we compute the visibility and expected termination depth at each location and in either direction along each ray for use as supervision for the visibility MLP. We minimize the sum of three losses:

$$L = \frac{1}{|R|} \sum_{r \in R} (\tilde{L}(r) - L(r))^2 + \frac{1}{|R|} \sum_{r \in R} (\tilde{V}(r) - V(r))^2 + \frac{1}{|R|} \sum_{r \in R} (\tilde{D}(r) - D(r))^2 \quad (16)$$

where $\sigma(x) = x/(1+x)$ is a tone-mapping operator [13], $L(r)$ and $\tilde{L}(r)$ are the ground truth and predicted camera ray radiance values (ground-truth values are simply the colors of input image pixels), $\tilde{V}(r)$ and $\tilde{D}(r)$ are the predicted visibility and expected termination depth from our visibility MLP given its current weights, $V(r)$ and $D(r)$ are the estimates of visibility and termination depth implied by the shape MLP given its current weights, and $\lambda = 20$ is the weight of the loss terms encouraging the visibility MLP to be consistent with the shape MLP. Note that the visibility MLP is not supervised using any “ground truth” visibility or termination depth — it is only optimized to be consistent with the NeRV’s current estimate of scene geometry, by evaluating Equations 5 and 6 using the densities emitted by the shape MLP. We apply a “stop gradient” to V and D in the last two terms of the loss, so the shape MLP is

Train Illum.	Hotdogs							
	Single Point		Colorful + Point		Ambient + Point		OLAT	
	PSNR	MS-SSIM	PSNR	MS-SSIM	PSNR	MS-SSIM	PSNR	MS-SSIM
NLT [61]	—	—	—	—	—	—	23.57	0.851
NeRF+LE	19.96	0.868	17.88	0.758	20.72	0.869	—	—
NeRF+Env	19.94	0.863	19.17	0.824	20.56	0.864	—	—
Bi <i>et al.</i> [3]	23.74	0.862	22.09	0.799	20.94	0.754	—	—
NeRV, NVF	23.93	0.860	24.37	0.885	25.14	0.892	—	—
NeRV, Trace	23.76	0.863	24.24	0.886	25.06	0.892	—	—

Train Illum.	Lego							
	Single Point		Colorful + Point		Ambient + Point		OLAT	
	PSNR	MS-SSIM	PSNR	MS-SSIM	PSNR	MS-SSIM	PSNR	MS-SSIM
NLT [61]	—	—	—	—	—	—	24.10	0.936
NeRF+LE	21.42	0.874	21.74	0.890	20.33	0.860	—	—
NeRF+Env	21.13	0.855	20.27	0.878	20.24	0.852	—	—
Bi <i>et al.</i> [3]	22.89	0.897	22.83	0.890	18.10	0.783	—	—
NeRV, NVF	22.78	0.866	23.82	0.899	23.32	0.894	—	—
NeRV, Trace	23.16	0.883	24.18	0.925	23.79	0.923	—	—

Train Illum.	Armadillo							
	Single Point		Colorful + Point		Ambient + Point		OLAT	
	PSNR	MS-SSIM	PSNR	MS-SSIM	PSNR	MS-SSIM	PSNR	MS-SSIM
NLT [61]	—	—	—	—	—	—	21.62	0.900
NeRF+LE	20.35	0.881	18.76	0.863	17.35	0.859	—	—
NeRF+Env	19.60	0.874	17.89	0.863	17.28	0.851	—	—
Bi <i>et al.</i> [3]	22.35	0.894	21.06	0.892	19.93	0.842	—	—
NeRV, NVF	21.14	0.882	22.80	0.910	22.80	0.897	—	—
NeRV, Trace	22.14	0.897	23.02	0.921	22.81	0.895	—	—

Table 1: Quantitative relighting and view synthesis results. For every scene, we train each method on three datasets that contain images of the scene under different illumination conditions, and compare the metrics of all variants on the same testing dataset. Please refer to Section 4 for details.

not encouraged to degrade its own performance to better match the output from the visibility MLP. We implement our model in JAX [5], and optimize using Adam [21] with a learning rate that begins at 10^{-5} and decays exponentially to 10^{-6} over the course of optimization (the other Adam hyperparameters are default values: $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 10^{-8}$). Each model is trained for 1 million iterations using 128 TPU cores, which takes 1 day to converge.

4. Results

NeRV outperforms prior work, particularly in its ability to recover relightable scene representations from images observed under complex lighting. We urge the reader to view our supplementary video to appreciate NeRV’s relighting and view synthesis results. In Table 1 we show performance for rendering images from novel viewpoints with lighting conditions not observed during training. We evaluate two versions of NeRV: **NeRV with Neural Visibility Fields (NeRV, NVF)** and **NeRV with Test-time Tracing (NeRV, Trace)**. Both methods use the same training procedure as described above, and differ only in how evaluation is performed: “NeRV, NVF” uses the same visibility approximations used during training at test time, while “NeRV, Trace” uses brute-force tracing to estimate visibility to point light sources to render sharper shadows at test time. We compare against the following baselines:

Neural Light Transport [61] (NLT) requires an input proxy geometry (which we provide by running marching

!""#%&'()*+,-.:/0123456789:;<=>?@A B C D E F G H I J K L M N O P Q R S T U V W X Y Z [\] ^ _ ` a b c d e f g h i j k l m n o p q r s t u v w x y z { | } ~ ¡ ¢ £ ¤ ¥ ¦ § ¨ © ª « ¬ ® ¯ ° ± ² ³ ´ µ ¶ · ¸ ¹ º » ¼ ½ ¾ ¿ À Á Â Ã Ä Å Æ Ç È É Ê Ë Ì Í Î Ï Ñ Ò Ó Ô Õ Ö × Ø Ù Ú Û Ü Ý Þ ß à á â ã

Figure 5: Both (b) Bi *et al.* [3] and (c) NeRV recover high-quality relightable models when trained on images illuminated by a single point source. However, for more complex lighting such as “Ambient+Point”, (d) Bi *et al.* fails as its brute-force visibility computation is unable to simulate the surrounding ambient lighting during training. Their model minimizes training loss by making the scene transparent and is thus unable to render convincing images for the “single point light” (row 1) or “colorful set of points lights” (row 2) conditions. (e) Because NeRV correctly simulates light transport, its renderings more closely resemble (a) the ground truth.

!""#%&'()*+,-.:/0123456789:;<=>?@A B C D E F G H I J K L M N O P Q R S T U V W X Y Z [\] ^ _ ` a b c d e f g h i j k l m n o p q r s t u v w x y z { | } ~ ¡ ¢ £ ¤ ¥ ¦ § ¨ © ª « ¬ ® ¯ ° ± ² ³ ´ µ ¶ · ¸ ¹ º » ¼ ½ ¾ ¿ À Á Â Ã Ä Å Æ Ç È É Ê Ë Ì Í Î Ï Ñ Ò Ó Ô Õ Ö × Ø Ù Ú Û Ü Ý Þ ß à á â ã

Figure 6: Modeling appearance changes due to lighting with a latent code does not generalize to lighting conditions unlike those seen during training. Here we train (b, c) the two latent code baselines (d) and NeRV on the “Ambient+Point” dataset. The latent code models are unable to produce convincing renderings at test time, while NeRV trained on the same data renders high-quality images.

cubes [27] on NeRFs [32] trained from images of each scene rendered with fixed lighting), and trains a convolutional network defined in an object’s texture atlas space to perform simultaneous relighting and view synthesis. Though our method just requires images with known but unconstrained lighting conditions for training, NLT requires multi-view images captured “One-Light-at-a-Time” (OLAT), where each viewpoint is rendered multiple times, once per light source. See the supplemental material for qualitative comparisons.

NeRF + Learned Embedding (NeRF+LE) and **NeRF + Fixed Environment Embedding (NeRF+Env)** represent appearance variation due to changing lighting using latent variables. Both augment the original NeRF model with an additional input of a 64-dimensional latent code corresponding to the scene lighting condition. These approaches are similar to “NeRF in the Wild” [28], which also uses a latent code to describe appearance variation due to variable light-

ing. NeRF+LE uses a PointNet [39] encoder to embed the position and color of each light, and NeRF+Env simply uses a flattened environment map as the latent code.

Neural Reflectance Fields (Bi et al. [3]) uses a similar neural volumetric representation as NeRV, with the critical difference that brute-force raymarching is used to compute visibilities. This approach is therefore unable to consider illumination from sources other than a single point light during training. At test time, when time and memory constraints are less restrictive, it computes visibilities to all light sources.

We train each method (other than NLT) on nine datasets. Each consists of 150 images of a synthetic scene (“Hotdogs”, “Lego”, or “Armadillo”) illuminated by one of 3 lighting conditions: 1) “Point” contains a single white point light randomly sampled on a hemisphere above the scene for each frame, representing a laboratory setup similar to that of Bi et al. [3]. 2) “Colorful + Point” contains a randomly-sampled point light, as well as a set of 8 colorful point lights whose locations and colors are fixed across all images in the dataset. This represents a challenging scenario with multiple strong light sources that cast shadows and tint the scene. 3) “Ambient + Point” contains a randomly-sampled point light, as well as a dim grey environment map. This represents a challenging scenario where scene points are illuminated from all directions. We separately train each method on each of these nine datasets and measure performance on the corresponding scene’s test set, which consists of 150 images of the scene under novel lighting conditions (containing either one or eight point sources) not observed during training, rendered from novel viewpoints not observed during training.

4.1. Discussion

Our method outperforms all baselines in experiments that correspond to challenging complex lighting conditions, and matches the performance of prior work in experiments with simple lighting. As visualized in Figure 5, the method of Bi et al. performs comparably to ours in the case it is designed for: images illuminated by a single point source. However, their model’s performance degrades when it is trained on datasets that have complex lighting conditions (“Colorful+Point” and “Ambient+Point” experiments in Table 1), as its forward model is unable to simulate light from more than a single source during training. As visualized in Figure 6, our method thoroughly outperforms both latent code baselines, as they are unable to generalize to lighting conditions that are unlike those seen during training. Our method generally matches or outperforms the NLT baseline, which requires a controlled laboratory lighting setup and substantially more inputs than all other methods (the multi-view OLAT dataset we use to train NLT contains eight times as many images as our other datasets, and the original NLT paper [61] uses 150 OLAT images per viewpoint).

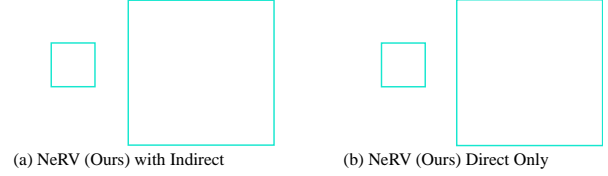


Figure 7: NeRV’s ability to simulate indirect illumination produces realistic details such as the additional brightness in the lego bulldozer’s cab due to interreflections.

(a) NeRV (Ours) with Analytic Normals (b) NeRV (Ours) with MLP-Predicted Normals

Figure 8: While obtaining surface normals (a) analytically or (b) as an output of the shape MLP produces similar renderings, analytic normals are much closer to the true surface normals.

Scene	Hotdogs		Lego		Armadillo	
	PSNR	MS-SSIM	PSNR	MS-SSIM	PSNR	MS-SSIM
Ours, No Indirect	24.43	0.861	23.06	0.888	21.27	0.878
Ours, MLP Normals	25.60	0.893	23.18	0.886	22.40	0.891
Ours, NVF	25.14	0.892	23.32	0.894	22.80	0.897
Ours, Trace	25.06	0.892	23.79	0.923	22.81	0.895

Table 2: Quantitative ablation results trained on the “Ambient+Point” dataset. Please refer to Section 4.2 for details.

4.2. Ablation Studies

We validate our choices of using analytic instead of MLP-predicted surface normals and simulating one-bounce indirect illumination through the ablation study reported in Table 2. We can see that modeling indirect illumination improves performance (Figure 7), even for our relatively simple scenes. Although using analytic instead of MLP-predicted normals is less numerically significant, Figure 8 shows that it results in more accurate estimated surface normals, which may be important for downstream graphics tasks.

5. Conclusion

We have demonstrated a method for recovering re-lightable neural volumetric representations from images of scenes illuminated by environmental and indirect lighting, by using a visibility MLP to approximate portions of the volume rendering integral that would otherwise be intractable to estimate during training by brute-force sampling. We believe that this work is an important initial foray into leveraging learned function approximation to alleviate the computational burden incurred by using rigorous physically-based differentiable rendering procedures for inverse rendering.

Acknowledgements We thank Peter Hedman for helpful feedback. MT is funded by an NSF Graduate Fellowship.

References

- [1] Jonathan T. Barron and Jitendra Malik. Shape, illumination, and reflectance from shading. *IEEE TPAMI*, 2015.
- [2] H. G. Barrow and J. M. Tenenbaum. Recovering intrinsic scene characteristics from images. *Computer Vision Systems*, 1978.
- [3] Sai Bi, Zexiang Xu, Pratul P. Srinivasan, Ben Mildenhall, Kalyan Sunkavalli, Miloš Hašan, Yannick Hold-Geoffroy, David Kriegman, and Ravi Ramamoorthi. Neural reflectance fields for appearance acquisition. *arXiv preprint arXiv:2008.03824*, 2020.
- [4] Sai Bi, Zexiang Xu, Kalyan Sunkavalli, Miloš Hašan, Yannick Hold-Geoffroy, David Kriegman, and Ravi Ramamoorthi. Deep reflectance volumes: Relightable reconstructions from multi-view photometric images. *ECCV*, 2020.
- [5] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, and Skye Wanderman-Milne. JAX: composable transformations of Python+NumPy programs, 2018.
- [6] Michael Bunnell. Dynamic ambient occlusion and indirect lighting. *NVIDIA GPU Gems 2*, 2004.
- [7] Zhang Chen, Anpei Chen, Guli Zhang, Chengyuan Wang, Yu Ji, Kiriakos N Kutulakos, and Jingyi Yu. A neural rendering framework for free-viewpoint relighting. *CVPR*, 2020.
- [8] Daniel Cohen and Zvi Sheffer. Proximity clouds—an acceleration technique for 3D grid traversal. *The Visual Computer*, 1994.
- [9] Boyang Deng, J. P. Lewis, Timothy Jeruzalski, Gerard Pons-Moll, Geoffrey Hinton, Mohammad Norouzi, and Andrea Tagliasacchi. Neural articulated shape approximation. *ECCV*, 2020.
- [10] Yue Dong, Guojun Chen, Pieter Peers, Jiawan Zhang, and Xin Tong. Appearance-from-motion: Recovering spatially varying surface reflectance under unknown lighting. *ACM TOG*, 2014.
- [11] Guan Gao, Guojun Chen, Yue Dong, Pieter Peers, Kun Xu, and Xin Tong. Deferred neural lighting: Free-viewpoint relighting from unstructured photographs. *ACM TOG*, 2020.
- [12] Kyle Genova, Forrester Cole, Aaron Sarna, Daniel Vlasic, William T. Freeman, and Thomas Funkhouser. Learning shape templates with structured implicit functions. *ICCV*, 2019.
- [13] Michaël Gharbi, Tzu-Mao Li, Miika Aittala, Jaakko Lehtinen, and Frédo Durand. Sample-based monte carlo denoising using a kernel-splatting network. *ACM TOG*, 2019.
- [14] Paul Green, Jan Kautz, and Frédo Durand. Efficient reflectance and visibility approximations for environment map rendering. *Eurographics*, 2007.
- [15] Romain Guy and Mathias Agopian. Physically based rendering in filament.
- [16] Richard Hartley and Andrew Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, 2nd edition, 2003.
- [17] Berthold K. P. Horn. Shape from shading: A method for obtaining the shape of a smooth opaque object from one view. Technical report, Massachusetts Institute of Technology, 1970.
- [18] Berthold K. P. Horn. Determining lightness from an image. *Computer Graphics and Image Processing*, 1974.
- [19] James T. Kajiya and Brian P. Von Herzen. Ray tracing volume densities. *SIGGRAPH*, 1984.
- [20] Brian Karis and Epic Games. Real shading in unreal engine 4. *Proc. Physically Based Shading Theory Practice*, 2013.
- [21] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015.
- [22] Edwin H. Land and John J. McCann. Lightness and retinex theory. *Journal of the Optical Society of America*, 1971.
- [23] 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. *CVPR*, 2020.
- [24] Zhengqi Li, Wenqi Xian, Abe Davis, and Noah Snavely. Crowdsampling the plenoptic function. *ECCV*, 2020.
- [25] Zhengqin Li, Zexiang Xu, Ravi Ramamoorthi, Kalyan Sunkavalli, and Manmohan Chandraker. Learning to reconstruct shape and spatially-varying reflectance from a single image. *ACM TOG*, 2018.
- [26] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *NeurIPS*, 2020.
- [27] William E. Lorensen and Harvey E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. *SIGGRAPH*, 1987.
- [28] Ricardo Martin-Brualla, Noha Radwan, Mehdi S.M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. NeRF in the wild: Neural radiance fields for unconstrained photo collections. *arXiv preprint arXiv:2008.02268*, 2020.
- [29] Nelson Max. Optical models for direct volume rendering. *IEEE TVCG*, 1995.
- [30] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3D reconstruction in function space. *CVPR*, 2019.
- [31] Moustafa Meshry, Dan B. Goldman, Sameh Khamis, Hugues Hoppe, Rohit Pandey, Noah Snavely, and Ricardo Martin-Brualla. Neural rerendering in the wild. *CVPR*, 2019.
- [32] 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. *ECCV*, 2020.
- [33] Giljoo Nam, Joo Ho Lee, Diego Gutierrez, and Min H Kim. Practical SVBRDF acquisition of 3D objects with unstructured flash photography. *ACM TOG*, 2018.
- [34] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3D representations without 3D supervision. *CVPR*, 2019.
- [35] Geoffrey Oxholm and Ko Nishino. Multiview shape and reflectance from natural illumination. *CVPR*, 2014.
- [36] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. *CVPR*, 2019.

- [37] Jeong Joon Park, Aleksander Holynski, and Steve Seitz. Seeing the world in a bag of chips. *CVPR*, 2020.
- [38] Matt Pharr, Wenzel Jakob, and Greg Humphreys. *Physically Based Rendering: From Theory to Implementation*. Morgan Kaufmann Publishers Inc., 3rd edition, 2016.
- [39] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. PointNet: Deep learning on point sets for 3D classification and segmentation. *CVPR*, 2017.
- [40] Ravi Ramamoorthi and Pat Hanrahan. A signal-processing framework for inverse rendering. *SIGGRAPH*, 2001.
- [41] Tobias Ritschel, Thomas Engelhardt, Thorsten Grosch, Hans-Peter Seidel, Jan Kautz, and Carsten Dachsbacher. Micro-rendering for scalable, parallel final gathering. *ACM TOG*, 2009.
- [42] Tobias Ritschel, Thorsten Grosch, Jan Kautz, and Stefan Müller. Interactive illumination with coherent shadow maps. *Eurographics Symposium on Rendering*, 2007.
- [43] Tobias Ritschel, Thorsten Grosch, Min H. Kim, Hans-Peter Seidel, Carsten Dachsbacher, and Jan Kautz. Imperfect shadow maps for efficient computation of indirect illumination. *ACM TOG*, 2008.
- [44] Hanan Samet. Implementing ray tracing with octrees and neighbor finding. *Computers & Graphics*, 1989.
- [45] Shen Sang and Manmohan Chandraker. Single-shot neural relighting and SVBRDF estimation. *ECCV*, 2020.
- [46] Yoichi Sato, Mark D. Wheeler, and Katsushi Ikeuchi. Object shape and reflectance modeling from observation. *SIGGRAPH*, 1997.
- [47] Carolin Schmitt, Simon Donne, Gernot Riegler, Vladlen Koltun, and Andreas Geiger. On joint estimation of pose, geometry and SVBRDF from a handheld scanner. *CVPR*, 2020.
- [48] Steven M. Seitz and Charles R. Dyer. Photorealistic scene reconstruction by voxel coloring. *IJCV*, 1999.
- [49] 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. *ICCV*, 2019.
- [50] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *NeurIPS*, 2020.
- [51] Vincent Sitzmann, Michael Zollhoefer, and Gordon Wetzstein. Scene representation networks: Continuous 3D-structure-aware neural scene representations. *NeurIPS*, 2019.
- [52] Peter-Pike Sloan, Jan Kautz, and John Snyder. Precomputed radiance transfer for real-time rendering in dynamic, low-frequency lighting environments. *Computer Graphics and Interactive Techniques*, 2002.
- [53] Noah Snavely, Steven M. Seitz, and Richard Szeliski. Photo tourism: Exploring photo collections in 3D. *ACM TOG*, 2006.
- [54] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *NeurIPS*, 2020.
- [55] Bruce Walter, Stephen R. Marschner, Hongsong Li, and Kenneth E. Torrance. Microfacet models for refraction through rough surfaces. *Rendering Techniques*, 2007.
- [56] Xin Wei, Guojun Chen, Yue Dong, Stephen Lin, and Xin Tong. Object-based illumination estimation with rendering-aware neural networks. *ECCV*, 2020.
- [57] Michael Weinmann and Reinhard Klein. Advances in geometry and reflectance acquisition (course notes). *SIGGRAPH Asia Courses*, 2015.
- [58] Rui Xia, Yue Dong, Pieter Peers, and Xin Tong. Recovering shape and spatially-varying surface reflectance under unknown illumination. *ACM TOG*, 2016.
- [59] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Ronen Basri, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *NeurIPS*, 2020.
- [60] Ye Yu and William A. P. Smith. InverseRenderNet: Learning single image inverse rendering. *CVPR*, 2019.
- [61] Xiuming Zhang, Sean Fanello, Yun-Ta Tsai, Tiancheng Sun, Tianfan Xue, Rohit Pandey, Sergio Orts-Escolano, Philip Davidson, Christoph Rhemann, Paul Debevec, Jonathan T. Barron, Ravi Ramamoorthi, and William T. Freeman. Neural light transport for relighting and view synthesis. *ACM TOG*, 2020.