

Inverse Rendering for Complex Indoor Scenes: Shape, Spatially-Varying Lighting and SVBRDF from a Single Image

Zhengqin Li[†], Mohammad Shafiei[†], Ravi Ramamoorthi[†], Kalyan Sunkavalli[‡], Manmohan Chandraker[†]

[†]University of California, San Diego

[‡]Adobe Research

Abstract

We propose a deep inverse rendering framework for indoor scenes. From a single RGB image of an arbitrary indoor scene, we create a complete scene reconstruction, estimating shape, spatially-varying lighting, and spatially-varying, non-Lambertian surface reflectance. To train this network, we augment the SUNCG indoor scene dataset with real-world materials and render them with a fast, high-quality, physically-based GPU renderer to create a large-scale, photorealistic indoor dataset. Our inverse rendering network incorporates physical insights – including a spatially-varying spherical Gaussian lighting representation, a differentiable rendering layer to model scene appearance, a cascade structure to iteratively refine the predictions and a bilateral solver for refinement – allowing us to jointly reason about shape, lighting, and reflectance. Experiments show that our framework outperforms previous methods for estimating individual scene components, which also enables various novel applications for augmented reality, such as photorealistic object insertion and material editing. Code and data will be made publicly available.

1. Introduction

A long-standing problem in computer vision is to reconstruct a scene—including its shape, lighting, and material properties—from a single image. This is an ill-posed task: these scene factors interact in complex ways to form images and multiple combinations of these factors may produce the same image [3]. As a result, previous work has often focused on subsets of this problem—shape reconstruction, illumination estimation, intrinsic images, etc.—or on restricted settings—single objects or objects from a specific class.

Our goal is a solution to a more general problem: from a single RGB image of an arbitrary indoor scene captured under uncontrolled conditions, we seek to reconstruct geometry, spatially-varying surface reflectance, and spatially-varying lighting. This is a challenging setting: indoor scenes demonstrate the entire range of real-world appearance, including arbitrary geometry and layouts, localized light sources that

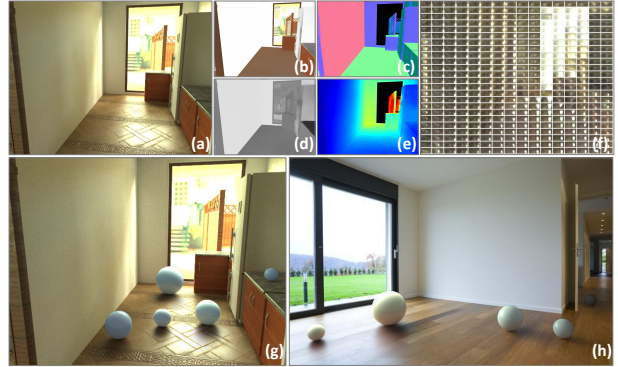


Figure 1. Given a single image of an indoor scene (a), we recover its diffuse albedo (b), normals (c), specular roughness (d), depth (e) and spatially-varying lighting (f). We train our model with high-quality synthetic images with photorealistic materials. By incorporating physical insights into our network structure, our predictions are of high enough equality to support applications like object insertion, even for specular objects (g) and in real images (h). Note the completely shadowed sphere on the extreme right.

lead to complex spatially-varying lighting effects, and complex, non-Lambertian surface reflectance. In this work we take a step towards providing a completely automatic, robust, holistic solution to this problem, thereby enabling a range of scene understanding and editing tasks. For example, in Figure 1(h), we use our scene reconstruction to enable photorealistic virtual object insertion; note how the inserted glossy spheres have realistic shading, shadowing caused by scene occlusions, and even reflections from the scene.

Driven by the success of deep learning methods on similar scene inference tasks (geometric reconstruction [19], lighting estimation [20], material recognition [11]), we propose training a deep convolutional neural network to regress these scene parameters from an input image. Ideally, the trained network should learn meaningful priors on these scene factors, and jointly model the interactions between them. In this work, we present two major contributions to address this.

Training deep neural networks requires large-scale, labeled training data. While datasets of real world geometry exist [17, 13], capturing real world lighting and surface



Figure 2. Comparison of object insertion results on a real image. Barron et al. [6] predict spatially varying log shading from a RGBD image, but their lighting representation does not preserve high frequency signal and cannot be used to render shadows and inter-reflections. Gardner et al. [20] predict a single lighting for the whole scene and therefore cannot model spatially varying indoor lighting. In contrast, our method solves the indoor scene inverse rendering problem in a holistic way, which results in photorealistic object insertion. The quality of our output may be visualized in a video for the example at the top, generated without any temporal constraints, at this [link](#).



Figure 3. A material editing example where we replace a material (on the surface of the kitchen counter-top) with a different one. Note the specular highlights on the surface, which can not be handled by conventional intrinsic decomposition methods since they do not recover the lighting direction. In contrast, we recover spatially-varying lighting and material properties.

reflectance at scale is non-trivial. Therefore, we use the SUNCG synthetic indoor scene dataset [53] that contains a large, diverse set of indoor scenes with complex geometry. However, the materials used in SUNCG are not realistic and the rendered images [60] are noisy. We address this by replacing SUNCG materials with high-quality, photorealistic SVBRDFs from a high-quality 3D material dataset [1]. We automatically map these SVBRDFs to SUNCG materials using deep features from a material estimation network, thus preserving scene semantics. We render the new scenes using a GPU-based global illumination renderer, to create high-quality input images. We also render the new scene reflectance and lighting and use it (along with the original geometry) to supervise our inverse rendering network.

An inverse rendering network would have to learn a model of image formation. The forward image formation model is well understood, and has been used in simple settings like

planar scenes and single objects [18, 38, 37, 40]. Indoor scenes are more complicated and exhibit challenging light transport effects like occlusions and inter-reflections. We address this by using a local lighting model—spatially-varying spherical gaussians (SVSGs). This bakes light transport effects directly into the lighting and makes rendering a purely local computation. We leverage this to design a fast, differentiable, *in-network* rendering layer that takes our geometry, SVBRDFs and SVSGs and computes radiance values. During training, we render our predictions and backpropagate the error through the rendering layer; this fixes the forward model, allowing the network to focus on the inverse task.

To the best of our knowledge, our work is the first demonstration of scene-level inverse rendering that truly accounts for complex geometry, lighting, materials, and light transport. Moreover, we demonstrate that we achieve results on par with state-of-the-art methods focused on specific tasks. For example, the diffuse albedo reconstructed using our method is competitive with a state-of-the-art intrinsic image method. Most importantly, by truly decomposing a scene into physically-based scene factors, we enable novel capabilities like photorealistic 3D object insertion and scene editing in images acquired in-the-wild. Figure 2 shows two object insertion examples on real indoor scene images. Since our method solves the inverse rendering problem in a holistic way, it achieves superior performances on object insertion compared with previous state-of-the-art methods [20, 6]. Figure 3 shows a material editing example, where we replace the material of a planar surface in a real image. Note that our method preserves spatially-varying specular highlights after changing the material. Such visual effects cannot be handled by traditional intrinsic decomposition methods.

2. Related Work

The problem of reconstructing shape, reflectance, and illumination from images has a long history in vision. It has been studied under different forms, such as intrinsic images (reflectance and shading from an image) [8] and shape-from-shading (shape, and sometimes reflectance, from an image) [25]. Here, we focus on *single* image methods.

Single objects. Many inverse rendering methods focus on reconstructing single objects. Even this problem is ill-posed and many methods assume some knowledge of the object in terms of known lighting [46, 27] or geometry [41, 49]. Other methods focus on specific object classes; for example, there are many methods that reconstruct facial shape, reflectance, and illumination using low-dimensional face models [12]. Recent methods have leveraged deep networks to reconstruct complex SVBRDFs from single images (captured under unknown environments) of simpler planar scenes [18, 37], objects of a specific class [40] or homogeneous BRDFs [43]. Other methods address illumination estimation [21]. We tackle the much harder case of large-scale scene modeling and do not assume scene information.

Barron and Malik [5] propose an optimization-based approach with hand-crafted priors to reconstruct shape, Lambertian reflectance, and distant illumination from a single image of an arbitrary object. Li et al. [38] tackle the same problem with a deep network and an object-specific rendering layer. Extending these methods to scenes is non-trivial because the light transport is significantly more complex.

Large-scale scenes. Previous work has looked at recognizing materials in indoor scenes [11] and decomposing indoor images into reflectance and shading layers [10, 36]. Techniques have also been proposed for single image geometric reconstruction [19] and lighting estimation [24, 20]. These methods estimate only one scene factor without modeling the rest of scene appearance, as we do.

Barron and Malik [6] reconstruct Lambertian reflectance and spatially-varying lighting but require an RGBD input image. Karsch et al. [32] propose a full-fledged scene reconstruction method that estimates geometry, Lambertian reflectance, and 3D lighting from a single image; however, they rely on extensive user input to annotate geometry and initialize lighting. Subsequently, they propose an automatic, rendering-based optimization method [33] that estimates all these scene factors. However, they rely on strong heuristics for their method that are often violated in the real world leading to errors in their estimates. In contrast, we propose a deep network that learns to predict geometry, complex SVBRDFs, and lighting in an end-to-end fashion.

Datasets. The success of deep networks has led to an interest in datasets for supervised training. This includes real world scans [17, 13] and synthetic shape [14] and scene [53]

datasets. All these datasets are either missing or have unrealistic material and lighting specifications. We build on the SUNCG dataset to improve its quality in this regard.

Differentiable rendering. A number of recent deep inverse rendering methods have incorporated in-network, differentiable rendering layers that are customized for simple settings: faces [51, 56], planar surfaces [18, 37], single objects [40, 38]. Some recent work has proposed differentiable general-purpose global illumination renderers [35, 15]; unlike our more specialized, fast rendering layer, these are too expensive to use for neural network training.

3. Dataset for Complex Indoor Scenes

A large-scale dataset is crucial for solving the complex task of inverse rendering in indoor scenes. It is extremely difficult, if at all possible, to acquire large-scale ground truth with spatially-varying material, lighting and global illumination. Thus, we resort to rendering a synthetic dataset, but must overcome significant challenges to ensure utility for handling real indoor scenes at test time. Existing datasets for indoor scenes are rendered with simpler assumptions on material and lighting. In this section, we describe our approach to photorealistically map our microfacet materials to SUNCG geometries [54], while preserving semantics. Further, rendering images with SVBRDF and global illumination, as well as ground truth for spatially-varying lighting, is computationally intensive, for which we design a custom GPU-accelerated renderer that outspeeds Mitsuba on a modern 16-core CPU by an order of magnitude.

3.1. Mapping photorealistic materials to SUNCG

Our goal is to map our materials to SUNCG geometries in a semantically meaningful way. The original materials in SUNCG dataset are represented by a Phong BRDF model [47] which is not suitable for complex materials [45]. Our materials, on the other hand, are represented by a physically motivated microfacet BRDF model [29], which consists of 1332 materials with high resolution 4096×4096 SVBRDF textures¹. This mapping problem is non-trivial: (i) specular lobes in SUNCG are not realistic [45, 55], (ii) an optimization-based fitting collapses due to local minima leading to serious over-fitting when used for learning and (iii) we must replace materials with similar semantic types while being consistent with geometry, for example, replace material on walls with other paints and on sofas with other fabrics. Thus, we devise a three-stage pipeline, summarized in Figure 4.

Step 1: Tileable texture synthesis Directly replacing SUNCG textures with our non-tileable ones will create artifacts near boundaries. Most frameworks for tileable texture

¹Please refer to Appendix D for details

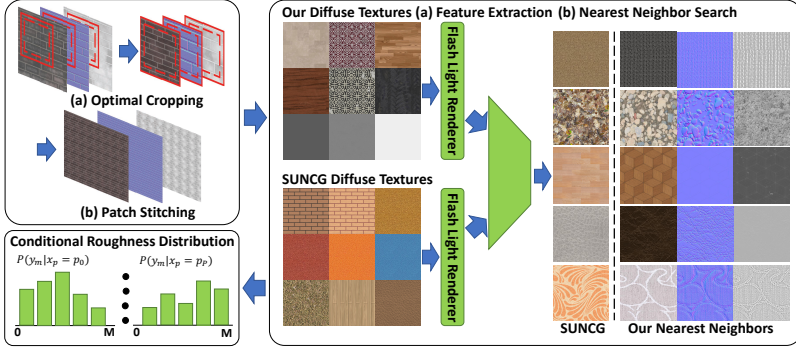


Figure 4. The pipeline of material mapping from original SUNCG materials with Phong BRDF to our microfacet BRDF. It has three steps. (Top left) Tileable texture synthesis to turns our SVBRDF textures into tileable ones. (Right) Spatially varying material mapping from SUNCG dataset with diffuse texture to our materials. (Bottom left) Homogeneous material mapping to convert specular parameters of homogeneous materials in SUNCG from Phong BRDF to our microfacet BRDF.

synthesis [39, 44] use randomized patch-based methods [4], which do not preserve structures such as sharp straight edges that are common for indoor scene materials such as bricks or wood floors. Instead, we first search for an optimal crop from our SVBRDF texture by minimizing gradients for diffuse albedo, normals and roughness perpendicular to the patch boundaries. We next find the best seam for tiling along the horizontal and vertical directions by modifying the graph cut method of [34] to encourage gradients to be similar at seams. Please refer to Appendix E for details on the energy design and examples of our texture synthesis.

Step 2: Mapping SVBRDFs Once our materials are tileable, we must use them to replace SUNCG ones in a semantically meaningful way. Since the specular reflectance of SUNCG materials is not realistic, we do this only for diffuse textures and directly use specularity from our dataset to render images. We manually divide 633 most common diffuse textures from SUNCG and from our entire dataset into 10 categories based on appearance and semantic labels, such as fabric, stone or wood. We render both sets of diffuse textures on a planar surface under flash light and use an encoder network similar to [37] to extract features, then use nearest neighbors to map the materials. We randomly choose from the 10 pre-computed nearest neighbors to render images in our dataset.

Step 3: Mapping homogeneous BRDFs To map homogeneous materials from SUNCG to ours, we keep the diffuse albedo unchanged and map specular Phong parameters to our microfacet model. Since the two lobes are very different, a direct fitting does not work. Instead, we compute a distribution of microfacet parameters conditioned on Phong parameters based on the mapping of diffuse textures, then randomly sample from that distribution to map specular parameters. Specifically, let $\mathbf{x}_P \in \mathcal{P}$ be specular parameters of Phong model and $\mathbf{y}_M \in \mathcal{M}$ be those of our microfacet BRDF. If a material from SUNCG has specular parameters $\mathbf{x}_P = \mathbf{p}_b$, we count the number of pixels in its 10 nearest neighbors from our dataset whose specular parameters are $\mathbf{y}_M = \mathbf{m}_a$. We sum up the number across the whole



Figure 5. The first row shows images rendered with materials from our dataset. The second and third rows are images rendered with the original materials from SUNCG dataset using Lambertian and Phong models. Images rendered with our materials have realistic specular highlights.

SUNCG dataset as $N(\mathbf{m}_a, \mathbf{p}_b)$. The probability of material with specular parameters \mathbf{y}_M given the original materials in SUNCG has specular parameters \mathbf{x}_P is defined as:

$$P(\mathbf{y}_M = \mathbf{m}_a | \mathbf{x}_P = \mathbf{p}_b) = \frac{N(\mathbf{p}_b, \mathbf{m}_a)}{\sum_{\mathbf{m}_c \in \mathcal{M}} N(\mathbf{p}_b, \mathbf{m}_c)}.$$

We sample the distribution as a piece-wise constant function and interpolate uniformly inside each bin to get continuous specular parameters of microfacet BRDF.

Comparative results Figure 5 shows a few scenes rendered with Lambertian, SUNCG Phong and our BRDF models. Images rendered with Lambertian BRDF do not have any specularity, those with Phong BRDF have strong but flat specular highlights, while ours are clearly more realistic. All the materials in our rendering are perfectly tiled and assigned to the correct objects, which demonstrates the effectiveness of our material mapping pipeline.

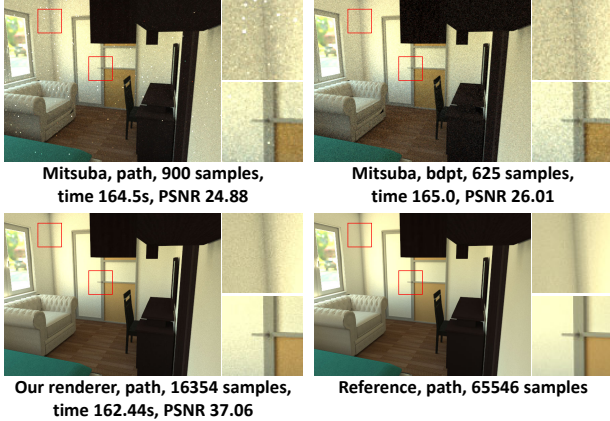


Figure 6. Comparisons of images rendered with Mitsuba and our GPU renderer in the same amount of time using path tracing. The quality of the image rendered by our renderer in less than three minutes is much better. It takes about 50 minutes for Mitsuba to achieve similar results.

3.2. Spatially Varying Lighting

To enable tasks such as object insertion or material editing, we must estimate lighting at every spatial location that encodes complex global interactions. Thus, our dataset must also render such ground truth. We do so by rendering a 16×32 environment map at the corresponding 3D point on object surfaces at every pixel.

In Figure 8, we show that an image obtained by integrating the product of this lighting and BRDF over the hemisphere looks very realistic, with high frequency specular highlights being correctly rendered. Note that global illumination and occlusion have already been baked into per pixel lighting, which makes it possible for a model trained on our lighting dataset to reason about those complex effects.

Enhancing lighting variations The SUNCG dataset is rendered with only one outdoor environment map and two area light intensities (400 for light bulbs and 0.5 for light shades). We add variations to the lighting to ensure generalizability, using 218 HDR outdoor panoramas from [23] and random RGB intensities for area lights.

3.3. Fast Physically-Based Rendering

To render high quality images with realistic appearances, it is necessary to use a physically based renderer that models complex light transport effects such as global illumination and soft shadows. However, current open source CPU renderers are too slow for creating a large dataset, especially to render per-pixel lighting. Thus, we implement our own physically-based GPU renderer using Nvidia OptiX [2]. To render a 480×640 image with 16384 samples per pixel, our renderer on Tesla V100 GPU needs 3-6 minutes, while Mitsuba on 16 cores of Intel i7-6900K CPU needs around 1

hour. Figure 6 compares images rendered with Mitsuba [26] and with our renderer using the same amount of time.

Rendered dataset We render 78794 HDR images, with 72220 used for training and 6574 for testing. The resolution of each image is 480×640 . We also render per pixel ground-truth lighting for 26719 images in the training set and all images in the testing set, at a spatial resolution of 120×160 . Our dataset and renderer will be made publicly available.

4. Network Design

Estimating material, geometry and lighting from a single indoor image is an extremely ill-posed problem, which we solve using priors learned by our physically-motivated deep network (architecture shown in Figure 7). Our network consists of cascaded stages of a SVBRDF and geometry predictor, a spatially-varying lighting predictor and a differentiable rendering layer, followed by a bilateral solver for refinement.

Material and geometry prediction The input to our network is a single gamma-corrected low dynamic range image I , stacked with a predicted three-channel segmentation mask $\{\tilde{M}_o, \tilde{M}_a, \tilde{M}_e\}$ that separates pixels of object, area lights and environment map, where $(\tilde{\cdot})$ represents predictions. The mask is obtained through a pre-trained network and useful since some predictions are not defined everywhere (for example, BRDF is not defined on light sources). Inspired by [37, 38], we use a single encoder to capture correlations between material and shape parameters, obtained using four decoders for diffuse albedo (A), roughness (R), normal (N) and depth (D). Skip links are used for preserving details. Then the initial estimates of material and geometry are given by

$$\tilde{A}, \tilde{N}, \tilde{R}, \tilde{D} = \text{MGNet}_0(I, M). \quad (1)$$

Spatially Varying Lighting Prediction Inverse rendering for indoor scenes requires predicting spatially varying lighting for every pixel in the image. Using an environment map as the lighting representation leads to a very high dimensional output space, that causes memory issues and unstable training due to small batch sizes. Spherical harmonics are a compact lighting representation that have been used in recent works [28, 38], but do not efficiently recover high frequency lighting necessary to handle specular effects [48, 9]. Instead, we follow pre-computed radiance transfer methods [57, 22, 59] and use isotropic spherical Gaussians that approximate all-frequency lighting with a smaller number of parameters. We model the lighting as a spherical function $L(\eta)$ approximated by the sum of spherical Gaussian lobes:

$$L(\eta) = \sum_{k=1}^K F_k G(\eta; \xi_k, \lambda_k), \quad G(\eta; \xi, \lambda) = e^{\lambda(1-\eta \cdot \xi)}, \quad (2)$$

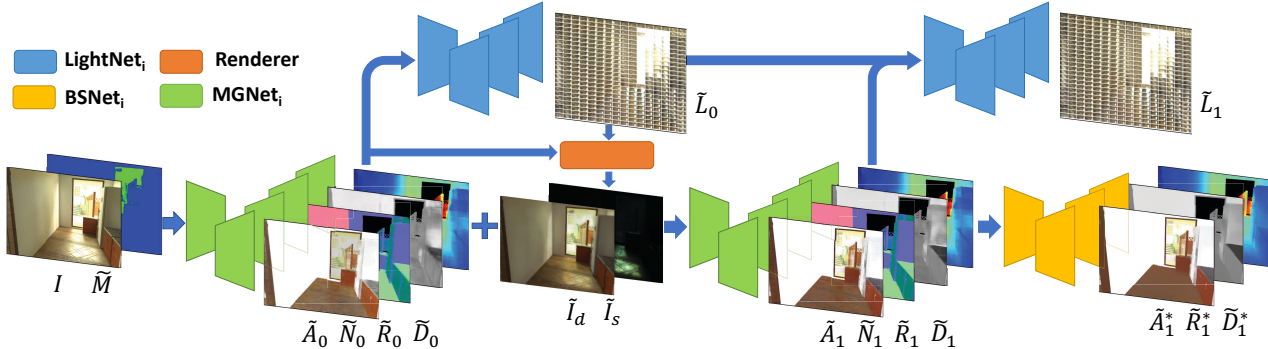
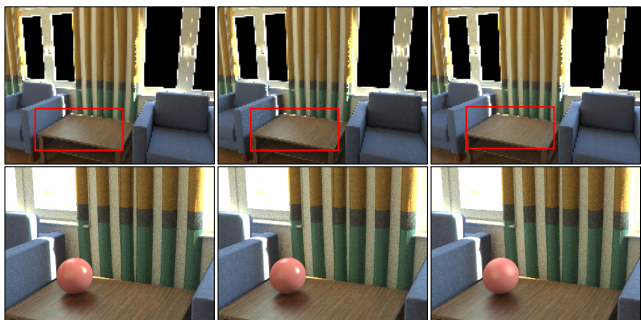


Figure 7. Our network design consists of a cascade, with one encoder-decoder for material and geometry prediction and another one for spatially-varying lighting, along with a physically-based differentiable rendering layer and a bilateral solver for refinement.



16x32x3 Environment map (1536 parameters) 12 spherical Gaussian lobes (72 parameters) 4 order spherical harmonic (75 parameters)

Figure 8. Comparisons of images rendered with lighting approximations. The first row: images rendered by our rendering layer using ground-truth normals and materials but with different lighting representations. The second row: inserting a sphere into the scene. In both examples, we can clearly see that spherical Gaussians can recover high frequency lighting much better with fewer parameters.

where η and ξ are vectors on the unit sphere \mathcal{S}^2 , F_k controls RGB color intensity and λ controls the bandwidth.

Each spherical Gaussian lobe is represented by 6 parameters $\{\xi_k, \lambda_k, F_k\}$. Figure 8 compares the images rendered with a 12-spherical Gaussian lobes approximation (72 parameters) and a fourth-order spherical harmonics approximation (75 parameters). Quantitative comparisons of lighting approximation and rendering errors are in Table 1. It is evident that even using fewer parameters, the spherical Gaussian lighting performs better, especially close to specular regions.

Our novel lighting prediction network, $\mathbf{LightNet}_0(\cdot)$, accepts predicted material and geometry as input, along with the image. It uses a shared encoder and separate decoders with a \tanh layer to predict:

$$\{\bar{\xi}_k\}, \{\bar{\lambda}_k\}, \{\bar{F}_k\} = \mathbf{LightNet}_0(I, \tilde{M}, \tilde{A}, \tilde{N}, \tilde{R}, \tilde{D}). \quad (3)$$

These original low dynamic range parameters are mapped to high dynamic range parameters $\{\bar{\xi}_k\}$, $\{\bar{\lambda}_k\}$ and $\{\bar{F}_k\}$

	Lighting ($\log L_2$)	Image (L_2)
SH (75 para.)	4.43	8.6×10^{-3}
SG (72 para.)	1.56	7.6×10^{-3}

Table 1. Quantitative comparison of using spherical harmonic (SH) and spherical Gaussian (SG) for lighting representation. From left to the right, the average error when using each representation to approximate per pixel lighting in Figure 8, the MSE of the rendered images. Again, spherical Gaussian performs better.

through following non-linear transformation.

$$\tilde{\xi}_k = \frac{\bar{\xi}_k}{\|\bar{\xi}_k\|_2} \quad (4)$$

$$\tilde{\lambda}_k = \tan\left(\frac{\pi}{4}(\bar{\lambda}_k + 1)\right) \quad (5)$$

$$\tilde{F}_k = \tan\left(\frac{\pi}{4}(\bar{F}_k + 1)\right). \quad (6)$$

Thus, our final predicted lighting is HDR, which is important for applications like relighting and material editing.

Differentiable rendering layer Our dataset in Section 3 provides ground truth for all scene components. But to model realistic indoor scene appearance, we additionally use a differentiable in-network rendering layer to mimic the image formation process, thereby weighting those components in a physically meaningful way. We implement this layer by numerically integrating the product of SVBRDF and spatially-varying lighting over the hemisphere. Let $l_{ij} = l(\phi_i, \theta_j)$ be a set of light directions sampled over the upper hemisphere, and v be the view direction. The rendering layer computes the diffuse image \tilde{I}_d and specular image \tilde{I}_s as:

$$\tilde{I}_d = \sum_{i,j} f_d(v, l_{ij}; \tilde{A}, \tilde{N}) L(l_{ij}; \{\xi_k, \lambda_k, F_k\}) \cos \theta_j d\omega, \quad (7)$$

$$\tilde{I}_s = \sum_{i,j} f_s(v, l_{ij}; \tilde{R}, \tilde{N}) L(l_{ij}; \{\xi_k, \lambda_k, F_k\}) \cos \theta_j d\omega, \quad (8)$$

where $d\omega$ is the differential solid angle. We sample 16×8 lighting directions. While this is relatively low resolution, we empirically find, as shown in Figure 8, that it is sufficient to recover most high frequency lighting effects.

Loss Functions Our loss functions incorporate physical insights. We first observe that two ambiguities are difficult to resolve: the ambiguity between color and light intensity, as well as the scale ambiguity of single image depth estimation. Thus, we allow the related loss functions to be scale invariant. For material and geometry, we use the scale invariant L_2 loss for diffuse albedo (\mathcal{L}_A), L_2 loss for normal (\mathcal{L}_N) and roughness (\mathcal{L}_R) and a scale invariant log-encoded loss for depth ($\mathcal{L}(D)$) due to its high dynamic range:

$$\mathcal{L}_D = \|(\log(D+1) - \log(c_d \tilde{D} + 1)) \odot (M_a + M_o)\|_2^2, \quad (9)$$

where c_d is a scale factor computed by least squares regression. For lighting estimation, we find supervising both the environment maps and spherical Gaussian parameters is important for preserving high frequency details. Thus, we compute ground-truth spherical Gaussian lobe parameters by approximating the ground-truth lighting using the LBFGS method². We use the same scale invariant log-encoded loss as (9) for weights ($\{\mathcal{L}_{F_k}\}$), bandwidth ($\{\mathcal{L}_{\lambda_k}\}$) and lighting ($\{\mathcal{L}_L\}$), with an L_2 loss for direction (L_{ξ_k}). We also add a scale invariant L_2 rendering loss:

$$\mathcal{L}_{ren} = \|(I - c_{diff} \tilde{I}_d - c_{spec} I_s) \odot M_o\|_2^2 \quad (10)$$

where \tilde{I}_d and \tilde{I}_s are rendered using (7) and (8), respectively, while c_{diff} and c_{spec} are positive scale factors computed using least square regression. The final loss function is a weighted summation of the proposed losses:

$$\begin{aligned} \mathcal{L} = & \alpha_A \mathcal{L}_A + \alpha_N \mathcal{L}_N + \alpha_R \mathcal{L}_R + \alpha_D \mathcal{L}_D + \alpha_L \mathcal{L}_L \\ & \alpha_{ren} \mathcal{L}_{ren} + \sum_{k=1}^K \alpha_\lambda \mathcal{L}_{\lambda_k} + \alpha_\xi \mathcal{L}_{\xi_k} + \alpha_F \mathcal{L}_{F_k}. \end{aligned} \quad (11)$$

Refinement using bilateral solver We use an end-to-end trainable bilateral solver to impose a smoothness prior [7, 36]. The inputs to a bilateral solver include the prediction, the estimated diffuse albedo \tilde{A} as a guidance image, and confidence map C . We train a shallow network with three sixteen-channel layers for confidence map predictions. Let $\mathbf{BS}_X(\cdot)$ be the bilateral solver and $\mathbf{BSNet}_X(\cdot)$ be the network for confidence map predictions where $X \in \{A, R, D\}$. We do not find refinement to have much effect on normals. The refinement process is:

$$\tilde{X} = \mathbf{BSNet}_X(\tilde{X}, I, \tilde{M}), \quad X \in \{A, R, D\} \quad (12)$$

$$\tilde{X}^* = \mathbf{BS}_X(\tilde{X}; C_X, \tilde{A}) \quad (13)$$

where we use (*) for predictions after refinement.

²Please refer to Appendix F for details on how we compute ground-truth spherical Gaussian parameters.

Cascade Network Akin to recent works on high resolution image synthesis [31, 16] and inverse rendering [38], we introduce a cascaded network that progressively increases resolution and iteratively refines the predictions through global reasoning. We achieve this by sending both the predictions and the rendering layer applied on the predictions to the next cascade stages, $\mathbf{MGNet}_1(\cdot)$ for material and geometry and $\mathbf{LightNet}_1(\cdot)$ for lighting, so that the network can reason about their differences.

$$\begin{aligned} \tilde{A}_1, \tilde{N}_1, \tilde{R}_1, \tilde{D}_1 = & \mathbf{MGNet}_1(I, \tilde{M}, \tilde{A}_0, \tilde{N}_0, \tilde{R}_0, \tilde{D}_0, \\ & c_{diff} \tilde{I}_d, c_{spec} \tilde{I}_s) \end{aligned} \quad (14)$$

$$\begin{aligned} \{\tilde{\xi}_k\}_1, \{\tilde{\lambda}_k\}_1, \{\tilde{F}_k\}_1 = & \mathbf{LightNet}_1(I, \tilde{M}, \tilde{A}_1, \tilde{N}_1, \tilde{R}_1, \\ & \tilde{D}_1, \{\tilde{\xi}_k\}_0, \{\tilde{\lambda}_k\}_0, \{\tilde{F}_k\}_0) \end{aligned} \quad (15)$$

Cascade stages have similar architectures as their initial network counterparts. One thing to notice is that we send low dynamic range lighting predictions $\{\xi_k\}_0, \{\lambda_k\}_0, \{F_k\}_0$ instead of the high dynamic range predictions, because we observe that it makes training more stable.

Training Details It is hard to train our whole pipeline end-to-end from scratch due to limited GPU memory, even with the use of group normalization [58]. So, we first train $\mathbf{MGNet}_i(\cdot)$ and $\mathbf{LightNet}_i(\cdot)$ separately with large batch sizes, fine-tune them together with smaller batch sizes, and finally train the bilateral solver. Please refer to Appendix G for training details and hyperparameter choices.

5. Experiments

Our experiments highlight the effectiveness of our dataset and network for single image inverse rendering in indoor scenes, through shape, material and lighting estimation. We achieve high accuracy on synthetic data and competitive performance on real images with respect to methods that focus only on a subset of those tasks. We conduct studies on the roles of various components in our pipeline. Finally, we illustrate applications such as high quality object insertion and material editing in real images that can only be enabled by our holistic solution to inverse rendering.

5.1. Analysis of Network and Training Choices

We study the effect of the cascade structure, joint training and refinement. Quantitative results for material and geometry predictions on the proposed dataset are summarized in Table 2, while those for lighting are shown in Table 3.

Cascade The cascade structure leads to clear gains for shape, BRDF and lighting estimation by iteratively improving and upsampling our predictions in Tables 2 and 3. This holds for both real data and synthetic data, as shown in Figure 9 and Figure 10. We observe that the cascade structure can effectively remove noise and preserve high frequency details for both materials and lighting. The errors in our

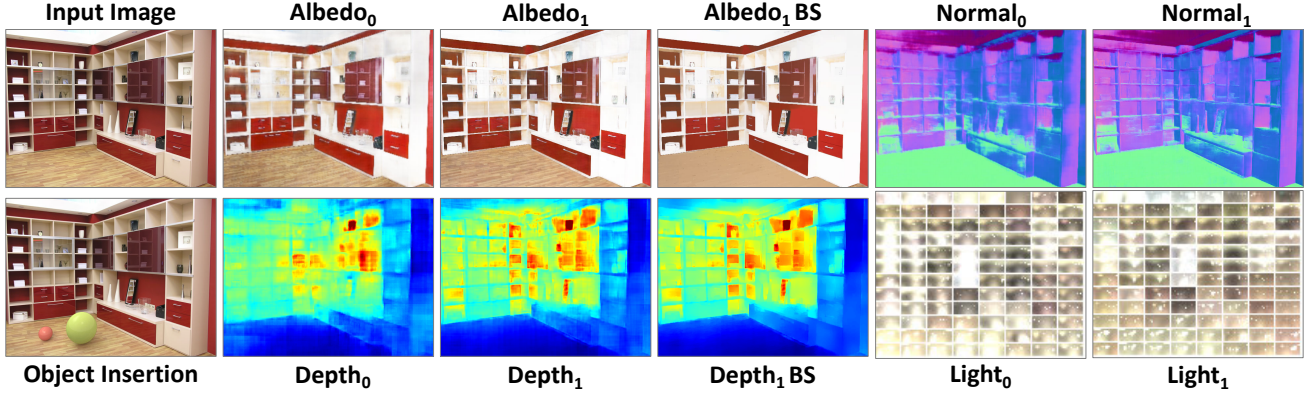


Figure 9. Impact of cascade and bilateral solver on a real example. Improvements are observed due to the cascade structure and bilateral solver. The estimates are accurate enough to insert a novel object with realistic global illumination effects.

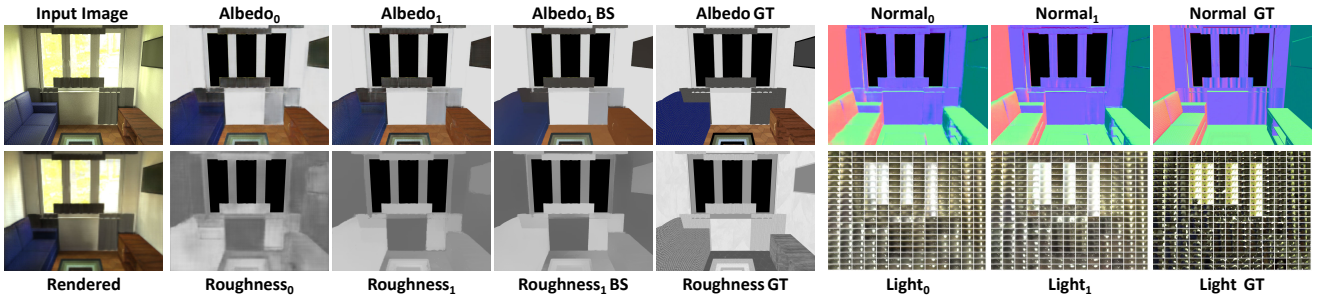


Figure 10. Impact of cascade and bilateral solver on a synthetic example. We observe that all the scene components benefit from cascaded estimation, while the bilateral solver is effective at refinement. The output of the rendering layer closely matches the input.

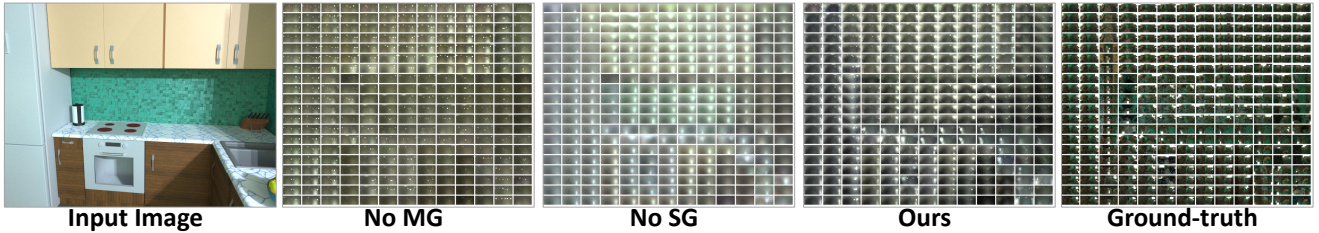


Figure 11. Comparison of lighting predictions. From left to the right are input image, **No MG**: without predicted material and geometry as input, **No SG**: without ground-truth spherical Gaussian parameters as supervision and our predictions and the ground-truth lighting.

shape, material and lighting estimates are low enough to photorealistically edit the scene to insert new objects, while preserving global illumination effects. In Figure 10, we observe that the image rendered using our predicted material, shape and lighting closely match the input image.

Joint training for inverse rendering Next we study whether BRDF, shape and lighting predictions can help improve each other. We compare jointly training the whole pipeline (“Joint”) using the loss in (11) and compare to independently training (“Ind”) each component $MGNet_i$ and $LightNet_i$. Quantitative errors on Tables 2 and 3 show that while errors for shape and BRDF prediction remain similar, those for rendering and lighting decrease. Next, we test lighting predictions without predicted BRDF as input for the

first level of cascade (“No MG”). Both quantitative results in Table 3 and qualitative comparison in Figure 11 demonstrate that the predicted BRDF and shape are important for the network to recover spatially varying lighting. We can see without the predicted material and geometry as input, the predicted lighting—especially the ambient color—does not sufficiently adapt spatially to the scene (possibly because of ambiguities between lighting and surface reflectance). This justifies our choice of jointly reasoning about shape, material and lighting. We also test lighting predictions with and without ground-truth SVSG parameters as supervision (“No SG”), finding that direct supervision leads to a sharper lighting prediction, which is shown in Figure 11.

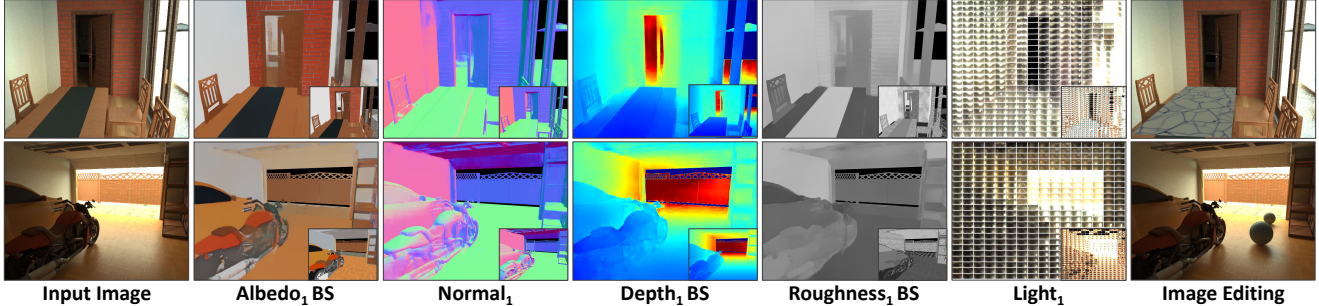


Figure 12. Results on our synthetic dataset. Given an input image, our estimated albedo, normals, depth, roughness and lighting are close to ground truth shown as insets. These are used for material editing (top) and object insertion (bottom).

	Cascade 0		Cascade 1		
	Ind.	Joint	Ind.	Joint	BS
$A(10^{-2})$	1.28	1.28	1.18	1.18	1.16
$N(10^{-2})$	4.91	4.91	4.91	4.51	4.51
$R(10^{-1})$	1.72	1.72	1.72	1.72	1.70
$D(10^{-2})$	8.06	8.00	7.29	7.26	7.20

Table 2. Quantitative comparisons of shape and material reconstructions on our test set. We use scale invariant L2 error for diffuse albedo (A), scale invariant \log^2 error for depth (D) and L2 error for normal (N) and roughness (R).

	Cascade 0				Cascade 1	
	No MG	No SG	Ind.	Joint	Ind.	Joint
L	2.83	2.85	2.54	2.50	2.49	2.43
$I(10^{-2})$	5.00	1.56	1.56	1.06	1.92	1.11

Table 3. Quantitative comparison of lighting predictions on test set. We use scale invariant L2 error for rendered image (I) and scale invariant \log^2 error for lighting (L).

Refinement Finally, we study the impact of the bilateral solver. Quantitative improvements over the second cascade stage in Table 2 are modest, which indicates that the network already learns good smoothness priors by that stage. This is shown in Figure 10, where the second level of cascade network generates smooth predictions for both material and lighting. But we find the qualitative impact of the bilateral solver to be noticeable on real images (for example, diffuse albedo in Figure 9), thus, we use it in all our real experiments.

Qualitative examples In Figure 12, we use a single input image from our synthetic test set to demonstrate depth, normal, SVBRDF and spatially-varying lighting estimation. The effectiveness is illustrated by low errors with respect to ground truth. Accurate shading and global illumination effects on an inserted object, as well as photorealistic editing of scene materials, show the utility of our decomposition.

5.2. Comparisons with Previous Works

We address the problem of holistic inverse rendering with spatially-varying material and lighting which has not been

Method	Training Set	WHDR
Ours (cascade 0)	Ours	23.29
Ours (cascade 1)	Ours	21.99
Ours (cascade 0)	Ours + IIW	16.83
Ours (cascade 1)	Ours + IIW	15.93
Li. et al[36]	CGI + IIW	17.5

Table 4. Intrinsic decomposition on the IIW dataset. Lower is better for the WHDR metric used here.

tackled earlier. Yet, it is instructive to compare our approach to prior ones that focus on specific sub-problems.

Intrinsic decomposition We compare two versions of our method on the IIW dataset [10] for intrinsic decomposition evaluation: our network trained on our data alone and our network fine-tuned on the IIW dataset. The results are tabulated in Table 4. We observe that the cascade structure is beneficial. We also observe a lower error compared to the prior work of [36], which indicates the benefit of our dataset that is rendered with a higher photorealism, as well as a network design that closely reflects physical image formation.

Lighting estimation We first compare to the method of Barron et al. [6] on our test set. Our scale-invariant shading errors on $\{R, G, B\}$ channels are $\{0.87, 0.86, 0.83\}$, compared to their $\{2.33, 2.10, 1.90\}$. Our shape, material and spatially-varying lighting estimation, together with a physically-motivated network trained on a realistic large-scale dataset, lead to this large improvement. Qualitative comparisons are shown in Figure 13, where we render specular spheres into the image using different lighting predictions. While our method can clearly capture the complex lighting variations and high frequency components, the spheres rendered with the predicted lighting of [6] are diffuse and have similar intensity across different regions of the image. Since only spherical harmonics parameters for log shading are predicted by [6], there is no physically correct way to turn its estimated spherical harmonics into environment lighting. Therefore, it cannot handle shadows and inter-reflections between object and the scene. Further, since only two orders of spherical harmonic parameters are predicted by [6], it cannot handle high frequency lighting.

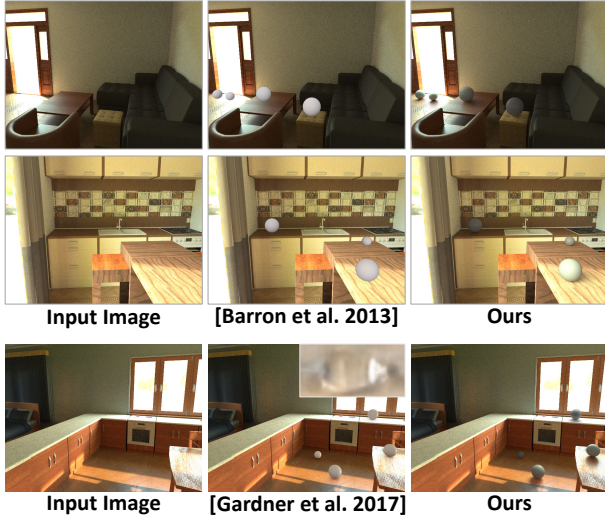


Figure 13. Comparison of object insertion with Barron et al. [6] (First two rows) and Gardner et al. [20] (third rows) on synthetic examples. We observe that rendered appearances from [6] are flat, while ours reflect the spatial variation in lighting at different parts of the scene. While method of [20] can preserve high frequency in the lighting, their lighting directions are not accurate and they could not handle spatially varying lighting.

Method	Mean($^{\circ}$)	Median($^{\circ}$)	Depth(Inv.)
Ours (cascade 0)	27.08	21.14	0.217
Ours (cascade 1)	26.33	20.21	0.206

Table 5. Normal and depth estimation on NYU dataset.

Next, we also compare with the work of Gardner et al. [20], which predicts a single environment lighting for the whole indoor scene. Quantitative results on our test set show that their mean $\log L_2$ error across the whole image is 3.34 while our $\log L_2$ error is 2.43. Qualitative results are shown in Figure 13. Since only one lighting for the whole scene is predicted by [20], no spatially-varying lighting effects can be observed.

In Figure 14, we compare our method with [5] and [20] on several real examples for object insertion in an image with spatially-varying illumination. It is clear that our method achieves a significant improvement in object insertion.

Depth and normal estimation We fine-tune our network, trained on our synthetic dataset, on NYU dataset as discussed in Appendix G. The test error on NYU dataset is summarized in Table 5. When testing depth error, we do not consider ground-truth depth value smaller than 1 or larger than 10 since they are outside the valid range of the sensor. When testing normal error, we mask out regions without accurate ground-truth normals. For both depth and normal prediction, the cascade structure consistently helps improve performance. Zhang et al. [60] achieve state-of-the-art performance for normal estimation using a more complex fine-

tuning strategy by choosing images with similar appearance as NYU dataset and with more than six times as much training data. Eigen et al. [19] achieve better results by using 120K frames of raw video data to train their network, while we pre-train on synthetic images with larger domain gap, using only use 795 images from NYU dataset for fine-tuning. Although we do not achieve competitive performance on this task, it’s not our main focus. Rather, we illustrate the wide utility of our proposed dataset and demonstrate estimation of factors of image formation good enough to support photorealistic augmented reality applications.

5.3. Novel Applications

Learning a disentangled shape, SVBRDF and spatially-varying lighting representation allows new applications that were hitherto not possible. We consider two of them here, object insertion and material editing. Before we discuss the two applications, we first describe how we resolve the ambiguity between scales of lighting and diffuse albedo.

Scales of lighting and diffuse albedo We use scale invariant loss for both diffuse albedo and lighting prediction. However, for real applications, we need to recover the scale of both diffuse albedo and lighting. Let c_a and c_l be the coefficients of diffuse albedo and lighting, respectively. Recall that our rendering layer outputs a diffuse image \tilde{I}_d and a specular image \tilde{I}_s . We can compute coefficients c_d and c_s to minimize the L_2 error between $c_d\tilde{I}_d + c_s\tilde{I}_s$ and input image I . Since our specular albedo is a constant, the scaling factor for our lighting prediction will be c_s and we have

$$c_l = c_s \quad (16)$$

$$c_a = \frac{c_d}{c_l} \quad (17)$$

However, for some images, specularity might be hard to observe, in which case we neglect the specularity term and simply compute the coefficient using

$$c_a = \frac{1}{\max(A)} \quad (18)$$

$$c_l = c_d/c_a. \quad (19)$$

That is, we set the scale of diffuse albedo c_a so that the largest albedo in the image is 1 and compute the coefficient of the lighting accordingly. To decide which strategy to use to compute the scale of lighting and albedo, we compute the following determinant when we regress c_d and c_s :

$$\mathcal{D} = \frac{(\tilde{I}_d \cdot \tilde{I}_d)(\tilde{I}_s \cdot \tilde{I}_s) - (\tilde{I}_d \cdot \tilde{I}_s)^2}{K}, \quad (20)$$

where K is the number of pixels in the image. If $\mathcal{D} > 1e^{-7}$, we use (16) and (17), otherwise we use (18) and (19) to compute the coefficient.

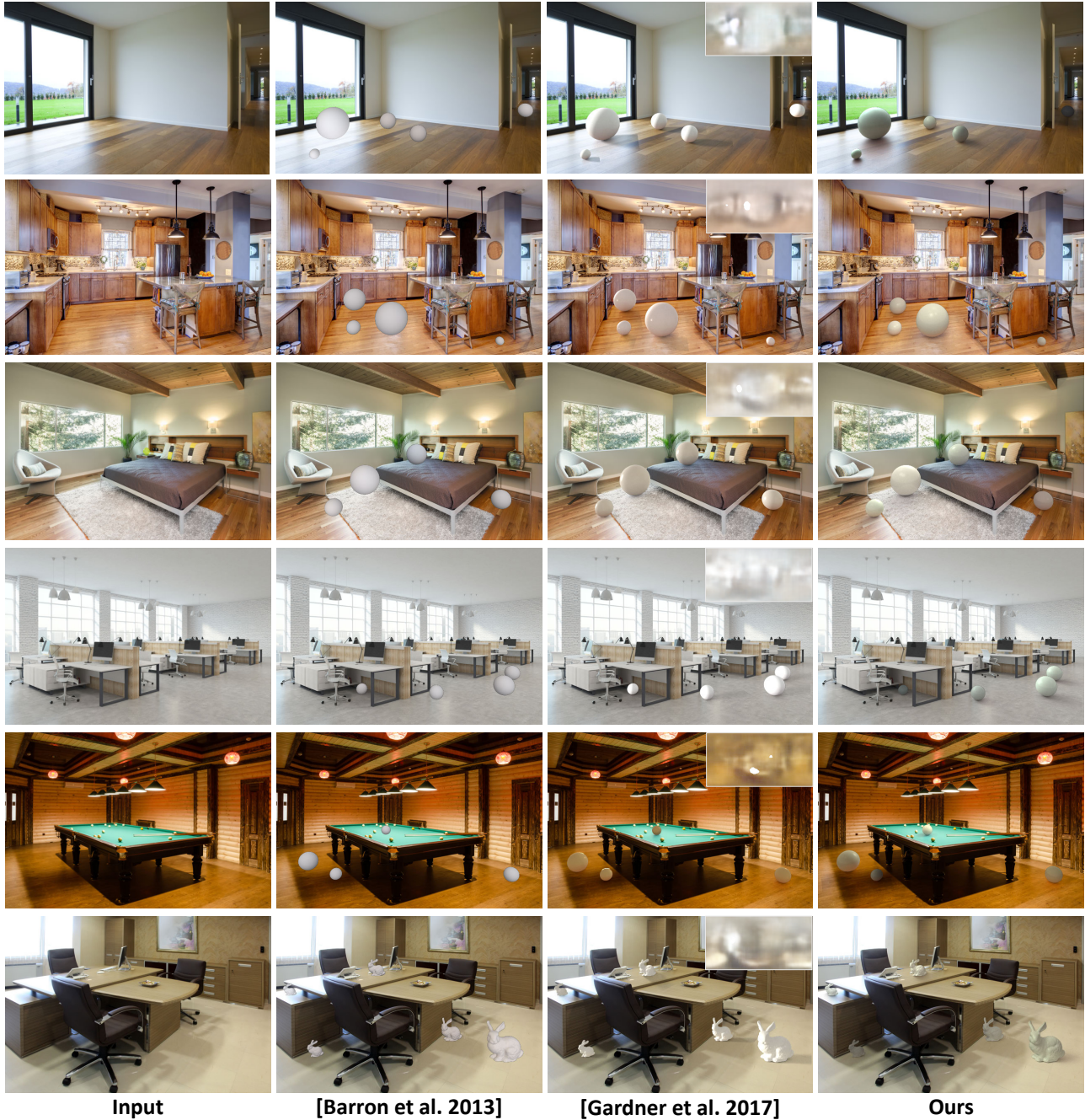


Figure 14. Object insertion examples and comparisons. Our proposed method estimates shape (depth and surface normals), spatially-varying complex reflectance (based on a micro-facet SVBRDF model) and spatially-varying lighting from a single image of an indoor scene. Given these estimates, we can insert virtual 3D objects into these images and produce photo-realistic results where the objects look like they truly belong in the scene. Note the shading and specular highlights on the inserted spheres, the realistic shadows cast on the ground, the reflections from the ground onto the spheres and adaptation of the appearance of the spheres to the local shadows and shading in the scene. We also compare with previous works of Barron et al. [6] and Gardner et al. [20] on real images. Note that in the results of [20], the shadows of some objects might be truncated by the plane we segment from the scene.

Object insertion To render a new object into the scene, we first crop a planar region and pick a point on that plane to place the new object. The orientation, diffuse albedo

and roughness value of the plane are all obtained from our predictions. We then render the plane and the object together using the lighting predicted at the point where we place the

object. We render the plane and the new object together to ensure inter-reflections between them are properly simulated. We compute a high resolution environment map (512×1024) from the estimated spherical Gaussian parameters so that even very glossy material can be correctly handled.

We render two images, I_{all} and I_{pl} and two binary masks, M_{obj} and M_{all} . I_{all} is the rendered image of plane and object and I_{pl} is the rendered image of the plane only. M_{obj} is the mask of the object and M_{all} is the mask covered both the cropped plane and the object. We then edit the region of object and the region of cropped plane separately. Let I be the original image and I_{new} be the new image with the new rendered object. For the object region, we directly use the intensities as rendered in I_{all} on the virtual object:

$$I_{new} \odot M_{obj} = I_{all} \odot M_{obj}. \quad (21)$$

For the remaining region on the plane, we blend in the original image intensities with the ratio of I_{all} and I_{pl} :

$$I_{new} \odot (M_{all} - M_{obj}) = I \odot \frac{I_{all}}{I_{pl}} \odot (M_{all} - M_{obj}). \quad (22)$$

All operations in the above relation are pixel-wise. This compositing procedure utilizes the idea of ratio (or quotient images) that has been used in the past for relighting [42, 50]. It ensures that global effects due to object-plane interaction, such as soft shadows, are visualized (since they are rendered in I_{all} but absent in I_{pl}), while keeping intensities consistent with the overall image. This suppresses high frequency artifacts in I_{all} that might be caused by minor errors in estimation of albedo, roughness and lighting, thereby achieving greater photorealism.

Figures 14, 2 and 1 show several examples of object insertion on real images. In all these examples, we render white glossy objects with diffuse albedo (0.8, 0.8, 0.8) and roughness value 0.2. We use white color so that the color of lighting and global illumination effects can be clearly observed. We keep the shape simple and the roughness value low to better demonstrate the high frequency component in the predicted lighting. To better demonstrate our performance, a video of an object moving around the scene rendered by our prediction can be found at this [link](#).

Material Editing Editing material properties of objects in a scene using a single photograph has applications for interior design and visualization. Our disentangled shape, material and lighting estimation allows rendering new appearances by replacing material and rendering using the estimated lighting. In Figure 3 and 15, we replace the material of a planar region with another kind of material and render the image using the predicted geometry and spatially varying lighting, where the spatially varying properties of the predicted lighting can be clearly observed. In the first example in Figure 3, we can see the specular high light in



Figure 15. Material editing. The left image is the original image and the right image is the rendered one with the material replaced in a part of the scene. We observe that the edited material looks photorealistic and even high frequency details from specular highlights and spatially-varying lighting are rendered well.

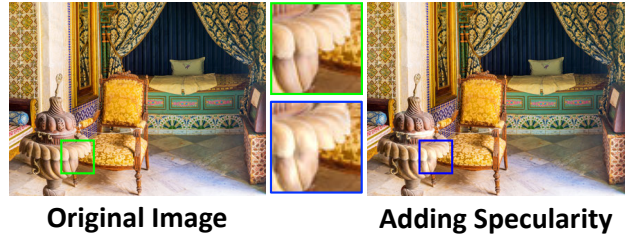


Figure 16. Changing the specularity of an object. We keep the predicted geometry, diffuse albedo and spatially varying lighting as predicted, but change the roughness value to 0.2 and re-render the object, leading to more prominent specular highlights.

the original image is preserved after changing the material, such specular high light effect can not be modeled by traditional intrinsic decomposition method since the direction of the incoming lighting is unknown. In Figure 16, we add specular highlight to the selected object by changing the roughness value to 0.2 and render the object with predicted diffuse albedo, geometry and spatially varying lighting. We compute the residual image before and after changing the roughness value and add it back to the original image. Even though the difference is quite subtle, we observe that the distribution of the specular highlight looks plausible.

6. Conclusion

We have presented the first holistic inverse rendering framework that estimates disentangled shape, SVBRDF and spatially-varying lighting, from a single image of an indoor scene. Insights from computer vision, graphics and deep convolutional networks are utilized to solve this challenging ill-posed problem. A GPU-accelerated renderer is used to synthesize a large-scale, realistic dataset with complex mate-

rials and global illumination. Our per-pixel SVSG lighting representation captures high frequency effects. Our network design imbibes intuitions such as a differentiable rendering layer, which are crucial for generalization to real images. Design choices such as a cascade structure and a bilateral solver lead to further benefits. Despite solving the joint problem, we obtain competitive results with respect to prior works that focus on constituent sub-problems, which highlights the impact of our dataset, representation choices and network design. We demonstrate object insertion and material editing applications on real images that capture global illumination effects, motivating applications in augmented reality and interior design.

A. Appendix Outline

We have presented a method to automatically disentangle a single image of an indoor scene into its constituent physical scene factors – geometry, spatially-varying reflectance, and illumination. In these appendices, we present more results, analyses and details. This includes: more challenging cases for our model (Appendix B and Appendix C), details about our SVBRDF model (Appendix D), dataset creation (Appendix E), our lighting model (Appendix F) and our network architecture and training details (Appendix G).

B. Generalization to Outdoor Scenes

In this section, we test how well our model, which is trained with synthetic indoor scenes only, generalizes to outdoor scenes. The qualitative results are shown in Figure 17. While the network tries to interpret the outdoor scene into a room surrounded by walls, we observe that the overall estimation of geometry, lighting and diffuse albedo look reasonable. We also try to insert a new object into the scene using our predictions following the pipeline proposed in Section 5, then compare with a state-of-the-art outdoor lighting estimation method [24]. As shown in the last two columns in Figure 17, the method of [24] can better preserve high-frequencies in outdoor illumination, which results in shadows with hard boundaries, while our method tends to predict more low frequency lighting. This is probably due to the domain gap between training and test images. However, we notice that our model can usually predict the direction of the incoming light correctly and the spatial variation in the lighting prediction of our method looks much more realistic compared to [24], which predicts a single lighting model for the whole image.

C. A Failure Case

While we observe largely successful object insertions in most experiments, some failure cases do occur. The ambiguity between albedo and lighting is a hard one to disentangle. In some cases, the albedo is estimated to be too bright

wall paint	stone wall	leather	stone floor	plastic
127	185	10	172	94
stone specular	ground	fabric	wood floor	wood
25	243	180	25	42

Table 6. The distribution of materials in our dataset, for the chosen semantic categories.

(dark), with the lighting correspondingly estimated as too dark (bright). An example is shown in Figure 18. Regardless, we emphasize that being able to estimate spatially-varying lighting along with SVBRDF and shape is an extremely hard problem, for which our network succeeds in an overwhelming number of experiments.

D. BRDF Model and Material Categories

Our microfacet model We use a physically motivated microfacet BRDF model in our dataset. Let A , N and R be the spatially-varying diffuse albedo, normal and roughness, respectively. The BRDF model $f(l, v; A, N, R)$ is:

$$f(l, v; A, N, R) = f_d(l, v; A, N) + f_s(l, v; N, R), \quad (23)$$

$$f_d(v, l; A, N) = \frac{A}{\pi}, \quad (24)$$

$$f_s(v, l; N, R) = \frac{D(h, R)F(v, h)G(l, v, N, R)}{4(N \cdot l)(N \cdot v)}, \quad (25)$$

where $f_d(\cdot)$ and $f_s(\cdot)$ are the diffuse and specular BRDF components. Here, v and l are view and lighting directions, and h is the half angle vector, while $D(h, R)$, $F(v, H)$ and $G(l, v, h, R)$ are the distribution, Fresnel and geometric terms respectively, defined as

$$\begin{aligned} D(h, R) &= \frac{\alpha^2}{\pi [(N \cdot h)^2(\alpha^2 - 1) + 1]^2} \\ \alpha &= R^2, \\ F(v, h) &= (1 - F_0)2^{-[5.55473(v \cdot h) + 6.8316](v \cdot h)}, \\ G(l, v, R, N) &= G_1(v, R, N)G_1(l, R, N), \\ G_1(l, R, N) &= \frac{N \cdot l}{(N \cdot l)(1 - k) + k}, \\ G_1(v, R, N) &= \frac{N \cdot v}{(N \cdot v)(1 - k) + k}, \\ k &= \frac{(R + 1)^2}{8}. \end{aligned}$$

We set $F_0 = 0.05$ as suggested in [30].

Material categories For mapping our materials on the SUNCG geometry in a manner consistent with semantics such as objects in the scene, we manually classified our dataset as well as the SUNCG materials into 10 categories. Samples from each dataset for these categories are shown in Figure 19. The number of material samples for each in our dataset are shown in Table 6.

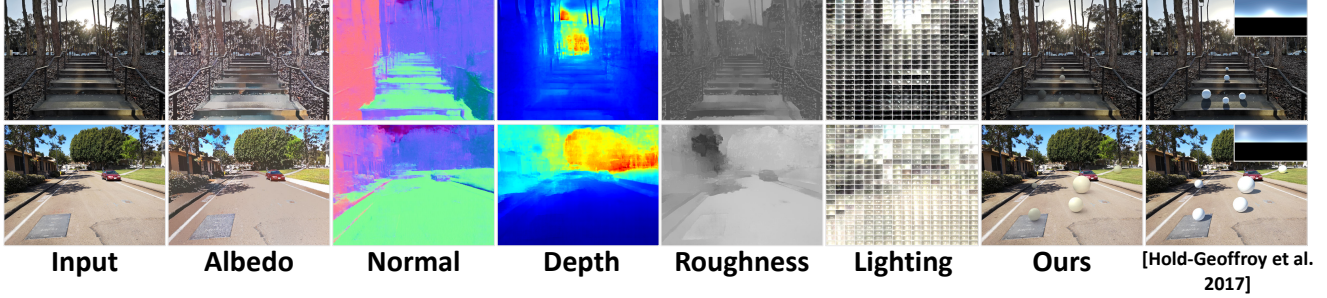


Figure 17. Our proposed method on outdoor scene. Even though our method is trained with synthetic indoor scenes only, it generalizes reasonably well to the outdoor scenes, which allows us to achieve reasonable object insertion results (the second last column) compared to the state-of-the-art [24] (the last column).

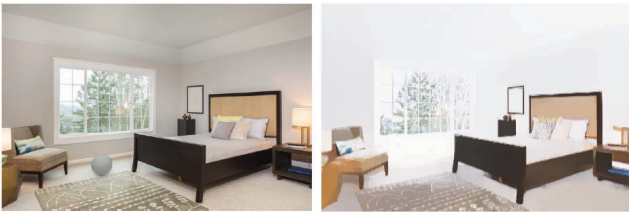


Figure 18. A failure case. (Left) The inserted object is rendered with an appearance that is darker than expected. (Right) The likely cause is an over-bright albedo estimation, which is traded off by a lighting estimate that has lower intensities.

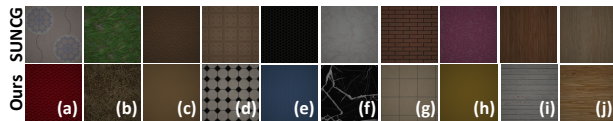


Figure 19. The ten material categories and the corresponding spatially varying diffuse textures from both SUNCG dataset and our dataset. From left to the right: (a) fabric, (b) ground, (c) leather, (d) stone floor, (e) plastic, (f) stone specular, (g) stone wall, (h) wall paint, (i) wood floor, (j) wood.

E. Tileable Texture Synthesis

We use graph-cut based approach to generate tileable texture, which has the advantages of keeping the original texture structures [34]. The overall process is summarized in Figure 20. We first crop smaller patch from the original SVBRDF textures and synthesize the cropped patch into tileable texture. Given the required size of the patch, we first globally search for the optimal patch by minimizing the gradient perpendicular to the boundary of the patches. More specifically, let \mathcal{B}_x and \mathcal{B}_y be the set of pixels on the horizontal and vertical boundaries of the patch. A_i , N_i and R_i are the diffuse color, normal and roughness at pixel i

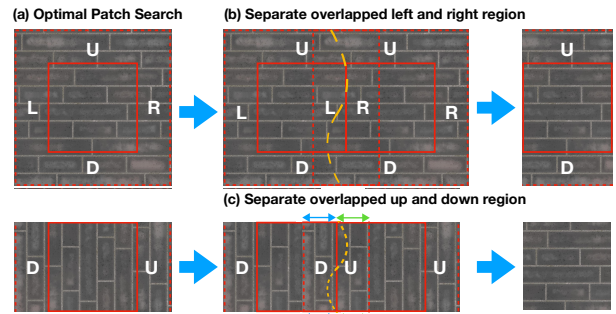


Figure 20. The pipeline of using graph-cut based method for tileable texture synthesis. We first stitch the left-right boundaries of textures and then the up-down boundaries. When stitching up-down boundaries, we add a hard constraint so that left-right boundaries will remain tileable.

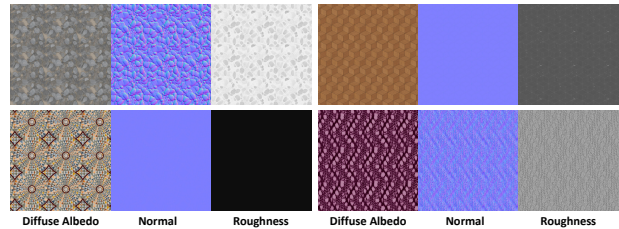


Figure 21. Results of tileable texture synthesis. Each image is generated by tiling 3×3 patches together.

respectively. We search for a patch so that

$$\sum_{i \in \mathcal{B}_y} \lambda_A \nabla_x A_i + \lambda_N \nabla_x N_i + \lambda_R \nabla_x R_i + \sum_{j \in \mathcal{B}_x} \lambda_A \nabla_y A_j + \lambda_N \nabla_y N_j + \lambda_R \nabla_y R_j \quad (26)$$

is minimized. The equation (26) can be efficiently computed using integral graph. The overall complexity of finding optimal patch will be $\mathcal{O}(K)$ where K is the number of pixels in the image. By minimizing (26), we avoid strong gradient near the boundaries of the patch, so that we can reduce

artifacts in tileable texture synthesis.

Once we find the patch, we crop not only the patch but also its surrounding regions. To make the patch tileable in x direction, we overlap the right and left surrounding regions and use graph-cut method to find the best seam to separate the overlapping regions by minimizing a customized energy function. Unlike energy function in [34] which encourages the value of pixels at seam to be similar to the value of pixels in the original textures, our energy function encourages the gradients of pixels at the seam to be similar to the gradients of pixels in the original textures. As in [34], we formulate the problem as a labeling problem. Let $I_{r,c}$ be an pixel in the overlapped texture map and $L_{r,c} \in \{1, 2\}$ be its label. With some abuse of notation, we define $I_{r,c}^i$ to be the value of pixel from patch i , $i \in \{1, 2\}$. The gradient across the patches i and j is defined as $\nabla_x I_{r,c}^{i,j} = I_{r,c+1}^j - I_{r,c}^i$. Then the loss $\mathcal{L}_I(L_{r,c} = 1, L_{r,c+1} = 2)$ is defined as

$$\min \left(\frac{\|\nabla_x I_{r,c}^{1,2} - \nabla_x I_{r,c}^{1,1}\|_1}{\max(\|\nabla_x I_{r,c}^{1,1}\|_1, 0.1)}, \frac{\|\nabla_x I_{r,c}^{1,2} - \nabla_x I_{r,c}^{2,2}\|_1}{\max(\|\nabla_x I_{r,c}^{2,2}\|_1, 0.1)} \right)$$

The final loss for is a weighted combination of losses from different texture map.

$$\lambda_A \mathcal{L}_A + \lambda_N \mathcal{L}_N + \lambda_R \mathcal{L}_R \quad (27)$$

To make the texture tileable in y direction, we repeat the above process by overlapping the up and down surrounding regions and finding the seam to separate them using graph-cut again. Notice that when separating the overlapping up and down regions, we need to make sure that the pixels at the right and left boundaries of the patches are from the same region so that the patch will remain tileable in x direction. We achieve this by adding an infinite smoothness term between every pair of pixels at the left and right boundaries in the same row so that they will always come from the same region. Figure 21 shows some texture synthesis examples. Each example is generated by tiling 3×3 original patches together. For each material from our dataset, we crop three patches of different sizes and the three patches will be considered as different materials in the following mapping SVBRDFs stage.

F. Ground Truth Spherical Gaussian Lobes

We compute ground-truth spherical Gaussian lobe parameters by approximating the environmental lighting using the LBFGS method. These parameters are used to supervise spatially varying lighting prediction. We use 12 lobes to approximate per pixel lighting. To facilitate the training process, we assign an order to the 12 lobes by constraining each lobe to be in the certain range of the hemisphere. We roughly divide the hemisphere into 2×6 regions. Following

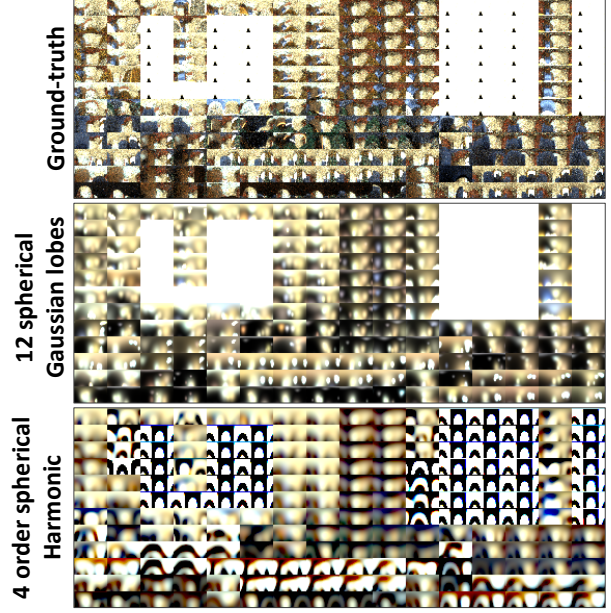


Figure 22. Comparison of approximating lighting with spherical harmonics and spherical Gaussian.

the notation in Section 4, we define $\{\theta_k\}$, $\{\phi_k\}$, $\{\lambda_k\}$ and $\{F_k\}$ to be the spherical Gaussian parameters where

$$\xi_k = (\sin \theta_k \cos \phi_k, \sin \theta_k \sin \phi_k, \cos \theta_k). \quad (28)$$

In order to add the constraints, we reparameterize the spherical Gaussian parameters with $\{\hat{\theta}_k\}$, $\{\hat{\phi}_k\}$, $\{\hat{\lambda}_k\}$ and $\{\hat{F}_k\}$ such that

$$\lambda_k = \exp(\hat{\lambda}_k), \quad (29)$$

$$F_k = \exp(\hat{F}_k), \quad (30)$$

$$\theta_k = a \tanh(\hat{\theta}_k) + b_k, \quad (31)$$

$$\phi_k = c \tanh(\hat{\phi}_k) + d_k, \quad (32)$$

where $a = \frac{3\pi}{8}$ and $c = \frac{\pi}{2}$ are scaling factors. Here, b_k and d_k are offset parameters that are computed as

$$b_k = \frac{\pi}{4}(k \bmod 2 + \frac{1}{2}), \quad (33)$$

$$d_k = \frac{\pi}{3}(k \bmod 6 + \frac{1}{2}) - \pi. \quad (34)$$

The initialization of the parameters are $\hat{\theta}_k = 0$, $\hat{\phi}_k = 0$, $\hat{F}_k = 0$ and $\hat{\lambda}_k = \log(\frac{\pi}{2})$. The loss function is the log-encoded loss as described in (9).

Figure 22 compares using spherical Gaussian and spherical harmonics to approximate the spatially varying lighting, which corresponds Figure 8 and Table 1 in Section 4. It is clearly observed that with a similar number of parameters, spherical Gaussians can better recover high frequency effects, resulting in a reconstructed spatially-varying lighting closer to ground truth.

G. Network Structures and Training Details

The network structures are summarized in Figure 23. Note that we use group normalization [58] instead of batch normalization so that we can train the network with smaller batch size. The padding size is dynamically assigned according to the feature map size so that the feature maps after up-sampling can be aligned with the feature maps coming from skip links. Therefore, our network can process image of arbitrary size without scaling and cropping. The network for spatially varying lighting predictions has more parameters because we find it necessary to achieve reasonable performances for this task.

We use Adam optimizer to train our network. Each level of cascade network is trained separately. To train cascade network of level i , we first train MGNet_i and LightNet_i separately and then fine-tune them together. The loss function to train MGNet_i is

$$\alpha_A \mathcal{L}_A + \alpha_N \mathcal{L}_N + \alpha_R \mathcal{L}_R + \alpha_D \mathcal{L}_D, \quad (35)$$

with various terms as defined in Sec. 4. We add the rendering loss \mathcal{L}_{ren} when training LightNet_i . The loss function to train LightNet_i is

$$\alpha_L \mathcal{L}_L + \alpha_{ren} \mathcal{L}_{ren} + \sum_{k=1}^K \alpha_{\lambda} \mathcal{L}_{\lambda_k} + \alpha_{\xi} \mathcal{L}_{\xi_k} + \alpha_F \mathcal{L}_{F_k}. \quad (36)$$

The loss function for fine-tuning the whole pipeline is defined in Eq. (11). Finally, the loss function to train BSNet is

$$\alpha_A \mathcal{L}_A + \alpha_R \mathcal{L}_R + \alpha_D \mathcal{L}_D. \quad (37)$$

All other hyper parameters including initial learning rate, training epochs and coefficients $\alpha_{(\cdot)}$ are summarized in Table 7 and Table 8. The learning rates are decreased by half every 10 epochs.

Fine-tuning on real datasets We use similar strategy to fine-tune on IIW dataset [10] and NYU dataset [52]. We take the trained model and fine-tune each level of cascade sequentially. The learning rate is $1e^{-4}$ and the batch size is 4 for the first level of cascade and 2 for the second level. In each iteration, we send two batches of images to the network, one from our synthetic dataset and the other from the real dataset. The MGNet_i and LightNet_i are trained in an end-to-end manner. The loss function for images from our dataset is the same as Eq. (11). The loss function for fine-tuning on NYU dataset is just a combination of the loss on each component. We add the rendering loss by comparing the rendered image with the input image:

$$\alpha_{ren} \mathcal{L}_{ren} + \alpha_N \mathcal{L}_N + \alpha_D \mathcal{L}_D. \quad (38)$$

When fine-tuning on IIW dataset, we include ordinal reflectance loss \mathcal{L}_{ord} which is the same as defined in [36]. The

loss function for images from IIW dataset is

$$\alpha_{ren} \mathcal{L}_{ren} + \alpha_{ord} \mathcal{L}_{ord}. \quad (39)$$

When training on NYU dataset, we also do data augmentation by randomly flipping, cropping and scaling the input images with a scale uniformly sampled from 0.8 to 1.2 since the dataset is relatively small.

References

- [1] Adobe Stock 3D assets. <https://stock.adobe.com/3d-assets>. 2
- [2] NVIDIA OptiX. <https://developer.nvidia.com/optix>. 5
- [3] E. H. Adelson and A. P. Pentland. Perception as bayesian inference. chapter The Perception of Shading and Reflectance, pages 409–423. 1996. 1
- [4] C. Barnes, E. Shechtman, A. Finkelstein, and D. B. Goldman. PatchMatch: A randomized correspondence algorithm for structural image editing. *ACM Transactions on Graphics (Proc. SIGGRAPH)*, 28(3), Aug. 2009. 3
- [5] J. Barron and J. Malik. Shape, illumination, and reflectance from shading. *PAMI*, 37(8):1670–1687, 2013. 3, 10
- [6] J. T. Barron and J. Malik. Intrinsic scene properties from a single rgb-d image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 17–24, 2013. 2, 3, 9, 10, 11
- [7] J. T. Barron and B. Poole. The fast bilateral solver. In *European Conference on Computer Vision*, pages 617–632. Springer, 2016. 7
- [8] H. G. Barrow and J. M. Tenenbaum. Recovering intrinsic scene characteristics from images. *Computer Vision Systems*, pages 3–26, 1978. 3
- [9] R. Basri and D. W. Jacobs. Lambertian reflectance and linear subspaces. *PAMI*, 25(2), 2003. 5
- [10] S. Bell, K. Bala, and N. Snavely. Intrinsic images in the wild. *ACM Transactions on Graphics (TOG)*, 33(4):159, 2014. 3, 9, 16
- [11] S. Bell, P. Upchurch, N. Snavely, and K. Bala. Material recognition in the wild with the materials in context database. *Computer Vision and Pattern Recognition (CVPR)*, 2015. 1, 3
- [12] V. Blanz and T. Vetter. A morphable model for the synthesis of 3d faces. In *SIGGRAPH*, pages 187–194, 1999. 3
- [13] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Niessner, M. Savva, S. Song, A. Zeng, and Y. Zhang. Matterport3D: Learning from RGB-D data in indoor environments. *International Conference on 3D Vision (3DV)*, 2017. 1, 3
- [14] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 3
- [15] C. Che, F. Luan, S. Zhao, K. Bala, and I. Gkioulekas. Inverse transport networks. *arXiv preprint arXiv:1809.10820*, 2018. 3
- [16] Q. Chen and V. Koltun. Photographic image synthesis with cascaded refinement networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1511–1520, 2017. 7

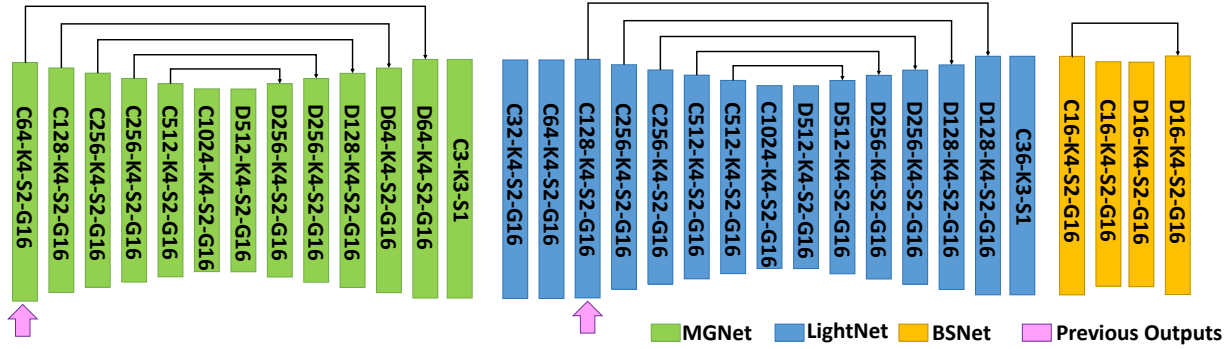


Figure 23. Our network architectures. $C/DX-KY-SZ-GW$ represents a convolution or transpose convolution layer with channel number X , kernel size Y , stride Z and group normalization with W channels in a group. **Previous Outputs** represents the place where we concatenate the predictions of previous level of cascade with current features.

	α_A	α_N	α_R	α_D	α_L	α_{ren}	α_λ	α_ξ	α_F	epochs	iters.	lr_{MG}	lr_{Light}	batch
MGNet₀	1.5	1.0	0.5	0.5	-	-	-	-	-	21	-	$1e^{-4}$	-	16
LightNet₀	-	-	-	-	10	10	$5e^{-4}$	1	0.5	15	-	-	$1e^{-4}$	4
Fine Tune ₀	7.5	5.0	2.5	2.5	10	10	$5e^{-4}$	1	0.5	-	4000	$1e^{-9}$	$1e^{-6}$	4
MGNet₁	1.5	1.0	0.5	0.5	-	-	-	-	-	8	-	$1e^{-4}$	-	6
LightNet₁	-	-	-	-	10	10	$5e^{-4}$	1	0.5	8	-	-	$1e^{-4}$	4
Fine Tune ₁	7.5	5.0	2.5	2.5	10	10	$5e^{-4}$	1	0.5	-	4000	$1e^{-9}$	$1e^{-6}$	3

Table 7. Hyper parameters for training **MGNet_i** and **LightNet_i**.

	α_A	α_R	α_D	iters.	lr_{BS}	batch
BSNet	1.5	0.5	0.5	600	$1e^{-4}$	6

Table 8. Hyper parameters for training **BSNet**.

- [17] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 2017. 1, 3
- [18] V. Deschaintre, M. Aittala, F. Durand, G. Drettakis, and A. Bousseau. Single-image svbrdf capture with a rendering-aware deep network. *ACM Transactions on Graphics (TOG)*, 37(4):128, 2018. 2, 3
- [19] D. Eigen and R. Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *ICCV*, 2015. 1, 3, 10
- [20] M.-A. Gardner, K. Sunkavalli, E. Yumer, X. Shen, E. Gambaretto, C. Gagné, and J.-F. Lalonde. Learning to predict indoor illumination from a single image. *ACM Trans. Graphics*, 9(4), 2017. 1, 2, 3, 10, 11
- [21] S. Georgoulis, K. Rematas, T. Ritschel, M. Fritz, T. Tuytelaars, and L. V. Gool. What is around the camera? In *ICCV*, 2017. 3
- [22] P. Green, J. Kautz, and F. Durand. Efficient reflectance and visibility approximations for environment map rendering. In *Computer Graphics Forum*, volume 26, pages 495–502. Wiley Online Library, 2007. 5
- [23] Y. Hold-Geoffroy, A. Athawale, and J.-F. Lalonde. Deep sky modeling for single image outdoor lighting estimation. In *CVPR*, 2019. 5
- [24] Y. Hold-Geoffroy, K. Sunkavalli, S. Hadap, E. Gambaretto, and J.-F. Lalonde. Deep outdoor illumination estimation. In *CVPR*, 2017. 3, 13, 14
- [25] B. K. P. Horn and M. J. Brooks, editors. *Shape from Shading*. MIT Press, Cambridge, MA, USA, 1989. 3
- [26] W. Jakob. Mitsuba renderer, 2010. <http://www.mitsuba-renderer.org>. 5
- [27] M. K. Johnson and E. H. Adelson. Shape estimation in natural illumination. In *CVPR*, 2011. 3
- [28] Y. Kanamori and Y. Endo. Relighting humans: occlusion-aware inverse rendering for fullbody human images. *SIG-GRAPH Asia*, 37(270):1–270, 2018. 5
- [29] B. Karis and E. Games. Real shading in unreal engine 4. *Proc. Physically Based Shading Theory Practice*, 4, 2013. 3
- [30] B. Karis and E. Games. Real shading in unreal engine 4. *Proc. Physically Based Shading Theory Practice*, 2013. 13
- [31] T. Karras, T. Aila, S. Laine, and J. Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017. 7
- [32] K. Karsch, V. Hedau, D. Forsyth, and D. Hoiem. Rendering synthetic objects into legacy photographs. *ACM Transactions on Graphics*, 30(6):1, 2011. 3
- [33] K. Karsch, K. Sunkavalli, S. Hadap, N. Carr, H. Jin, R. Fonte, M. Sittig, and D. Forsyth. Automatic scene inference for 3d object compositing. *ACM Transactions on Graphics*, (3):32:1–32:15, 2014. 3

- [34] V. Kwatra, A. Schödl, I. Essa, G. Turk, and A. Bobick. Graphcut textures: image and video synthesis using graph cuts. *TOG*, 22(3):277–286, 2003. 4, 14, 15
- [35] T.-M. Li, M. Aittala, F. Durand, and J. Lehtinen. Differentiable monte carlo ray tracing through edge sampling. *ACM Trans. Graph. (Proc. SIGGRAPH Asia)*, 37(6):222:1–222:11, 2018. 3
- [36] Z. Li and N. Snavely. Cgintrinsics: Better intrinsic image decomposition through physically-based rendering. In *ECCV*, pages 371–387, 2018. 3, 7, 9, 16
- [37] Z. Li, K. Sunkavalli, and M. Chandraker. Materials for masses: Svbrdf acquisition with a single mobile phone image. In *ECCV*, pages 72–87, 2018. 2, 3, 4, 5
- [38] Z. Li, Z. Xu, R. Ramamoorthi, K. Sunkavalli, and M. Chandraker. Learning to reconstruct shape and spatially-varying reflectance from a single image. In *SIGGRAPH Asia*, page 269. ACM, 2018. 2, 3, 5, 7
- [39] L. Liang, C. Liu, Y.-Q. Xu, B. Guo, and H.-Y. Shum. Real-time texture synthesis by patch-based sampling. *ACM Transactions on Graphics (ToG)*, 20(3):127–150, 2001. 3
- [40] G. Liu, D. Ceylan, E. Yumer, J. Yang, and J.-M. Lien. Material editing using a physically based rendering network. 2017. 2, 3
- [41] S. Lombardi and K. Nishino. Reflectance and natural illumination from a single image. In *ECCV*, 2012. 3
- [42] S. R. Marschner and D. P. Greenberg. Inverse lighting for photography. In *Color and Imaging Conference*, volume 1997, pages 262–265. Society for Imaging Science and Technology, 1997. 12
- [43] A. Meka, M. Maximov, M. Zollhofer, A. Chatterjee, H.-P. Seidel, C. Richardt, and C. Theobalt. Lime: Live intrinsic material estimation. In *CVPR*, 2018. 3
- [44] J. Moritz, S. James, T. S. Haines, T. Ritschel, and T. Weyrich. Texture stationarization: Turning photos into tileable textures. In *Computer Graphics Forum*, volume 36, pages 177–188. Wiley Online Library, 2017. 3
- [45] A. Ngan, F. Durand, and W. Matusik. Experimental analysis of brdf models. *Rendering Techniques*, 2005(16th):2, 2005. 3
- [46] G. Oxholm and K. Nishino. Shape and reflectance from natural illumination. In *ECCV*, 2012. 3
- [47] B. T. Phong. Illumination for computer generated pictures. *Communications of the ACM*, 18(6):311–317, 1975. 3
- [48] R. Ramamoorthi and P. Hanrahan. An efficient representation for irradiance environment maps. In *SIGGRAPH*, 2001. 5
- [49] F. Romeiro and T. Zickler. Blind reflectometry. In *ECCV*, 2010. 3
- [50] A. Shashua and T. Riklin-Raviv. The quotient image: Class-based re-rendering and recognition with varying illuminations. *PAMI*, 23(2):129–139, Feb. 2001. 12
- [51] Z. Shu, E. Yumer, S. Hadap, K. Sunkavalli, E. Shechtman, and D. Samaras. Neural face editing with intrinsic image disentangling. In *CVPR*, 2017. 3
- [52] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus. Indoor segmentation and support inference from rgb-d images. In *European Conference on Computer Vision*, pages 746–760. Springer, 2012. 16
- [53] S. Song, F. Yu, A. Zeng, A. X. Chang, M. Savva, and T. Funkhouser. Semantic scene completion from a single depth image. *Proceedings of 30th IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2, 3
- [54] S. Song, F. Yu, A. Zeng, A. X. Chang, M. Savva, and T. Funkhouser. Semantic scene completion from a single depth image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1746–1754, 2017. 3
- [55] T. Sun, H. W. Jensen, and R. Ramamoorthi. Connecting measured brdfs to analytic brdfs by data-driven diffuse-specular separation. *ACM Transactions on Graphics (TOG)*, 37(6):273, 2018. 3
- [56] A. Tewari, M. Zollhofer, H. Kim, P. Garrido, F. Bernard, P. Perez, and C. Theobalt. Mofa: Model-based deep convolutional face autoencoder for unsupervised monocular reconstruction. In *ICCV*, 2018. 3
- [57] Y.-T. Tsai and Z.-C. Shih. All-frequency precomputed radiance transfer using spherical radial basis functions and clustered tensor approximation. In *TOG*, volume 25, pages 967–976. ACM, 2006. 5
- [58] Y. Wu and K. He. Group normalization. In *ECCV*, pages 3–19, 2018. 7, 16
- [59] K. Xu, W.-L. Sun, Z. Dong, D.-Y. Zhao, R.-D. Wu, and S.-M. Hu. Anisotropic spherical gaussians. *ACM Transactions on Graphics (TOG)*, 32(6):209, 2013. 5
- [60] Y. Zhang, S. Song, E. Yumer, M. Savva, J.-Y. Lee, H. Jin, and T. Funkhouser. Physically-based rendering for indoor scene understanding using convolutional neural networks. *CVPR*, 2017. 2, 10