

This is a repository copy of *InverseRenderNet : Learning single image inverse rendering*.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/156112/>

Version: Accepted Version

Proceedings Paper:

Yu, Ye and Smith, William Alfred Peter orcid.org/0000-0002-6047-0413 (2020)

InverseRenderNet : Learning single image inverse rendering. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). 30th IEEE/CVF Conference on Computer Vision and Pattern Recognition 2019, 16-21 Jun 2019 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) . , USA , pp. 3150-3159.

<https://doi.org/10.1109/CVPR.2019.00327>

Reuse

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.

InverseRenderNet: Learning single image inverse rendering

Ye Yu and William A. P. Smith
Department of Computer Science, University of York, UK
{yy1571,william.smith}@york.ac.uk

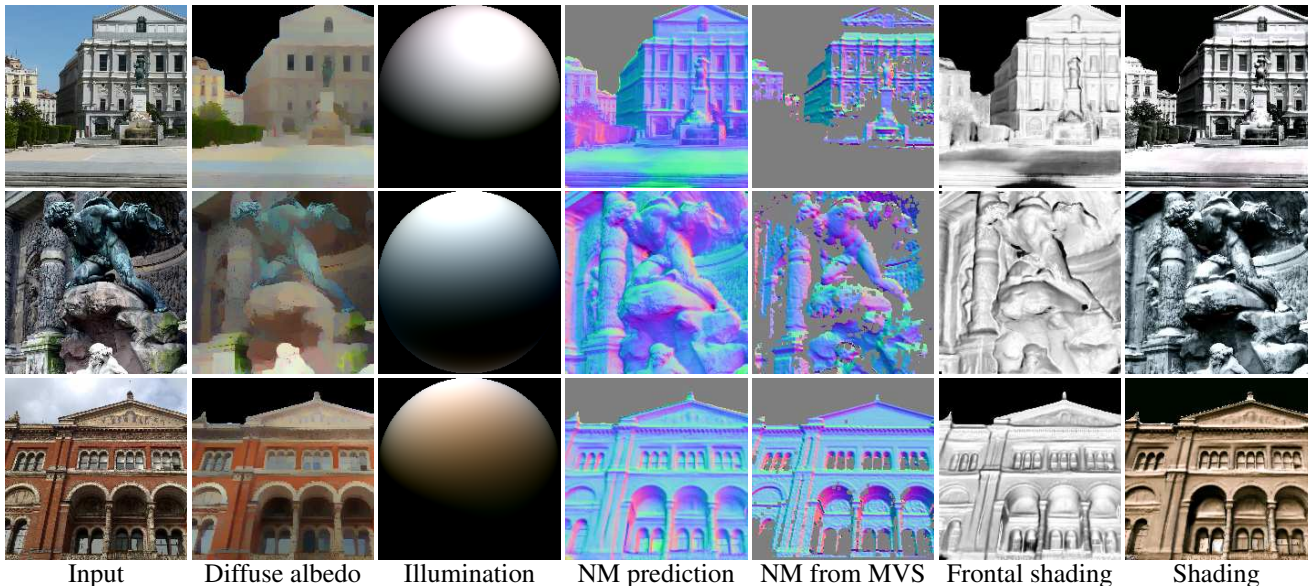


Figure 1: From a single image (col. 1), we estimate albedo and normal maps and illumination (col. 2-4); comparison multi-view stereo result from several hundred images (col. 5); re-rendering of our shape with frontal/estimated lighting (col. 6-7).

Abstract

We show how to train a fully convolutional neural network to perform inverse rendering from a single, uncontrolled image. The network takes an RGB image as input, regresses albedo and normal maps from which we compute lighting coefficients. Our network is trained using large uncontrolled image collections without ground truth. By incorporating a differentiable renderer, our network can learn from self-supervision. Since the problem is ill-posed we introduce additional supervision: 1. We learn a statistical natural illumination prior, 2. Our key insight is to perform offline multiview stereo (MVS) on images containing rich illumination variation. From the MVS pose and depth maps, we can cross project between overlapping views such that Siamese training can be used to ensure consistent estimation of photometric invariants. MVS depth also provides direct coarse supervision for normal map estimation. We believe this is the first attempt to use MVS supervision for learning inverse rendering.

1. Introduction

Inverse rendering is the problem of estimating one or more of illumination, reflectance properties and shape from observed appearance (i.e. one or more images). In this paper, we tackle the most challenging setting of this problem; we seek to estimate all three quantities from only a single, uncontrolled image. Specifically, we estimate a normal map, diffuse albedo map and spherical harmonic lighting coefficients. This subsumes two classical computer vision problems: (uncalibrated) shape-from-shading and intrinsic image decomposition.

Classical approaches [4, 29] cast these problems in terms of energy minimisation. Here, a data term measures the difference between the input image and the synthesised image that arises from the estimated quantities. We approach the problem as one of image to image translation and solve it using a deep, fully convolutional neural network. However, inverse rendering of uncontrolled, outdoor scenes is itself an unsolved problem and so labels for supervised learning

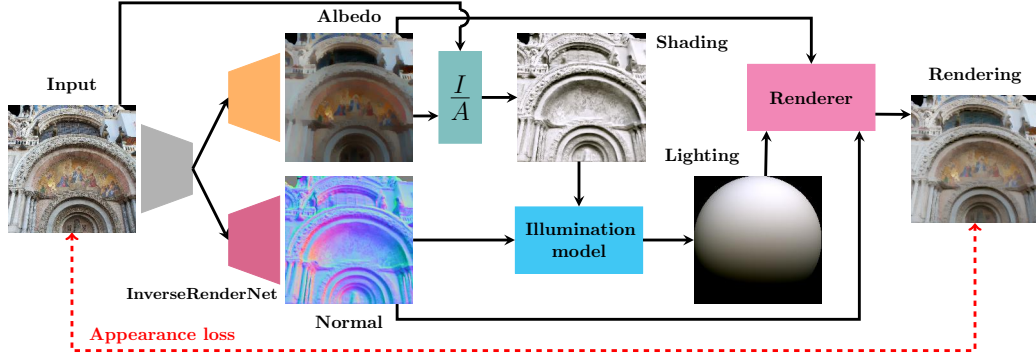


Figure 2: At inference time, our network regresses diffuse albedo and normal maps from a single, uncontrolled image and then computes least squares optimal spherical harmonic lighting coefficients. At training time, we introduce self-supervision via an appearance loss computed using a differentiable renderer and the estimated quantities.

are not available. Instead, we use the data term for self-supervision via a differentiable renderer (see Fig. 2).

Single image inverse rendering is an inherently ambiguous problem. For example, any image can be explained with zero data error by setting the albedo map equal to the image, the normal map to be planar and the illumination arbitrarily such that the shading is unity everywhere. Hence, the data term alone cannot be used to solve this problem. For this reason, classical methods augment the data term with generic [4] or object-class-specific [2] priors. Likewise, we also exploit priors during learning (specifically a statistical prior on lighting and a smoothness prior on diffuse albedo). However, our key insight that enables the CNN to learn good performance is to introduce additional supervision provided by an offline multiview reconstruction.

While photometric vision has largely been confined to restrictive lab settings, classical geometric methods are sufficiently robust to provide multiview 3D shape reconstructions from large, unstructured datasets containing very rich illumination variation [14, 17]. This is made possible by local image descriptors that are largely invariant to illumination. However, these methods recover only geometric information and any recovered texture map has illumination “baked in” and so is useless for relighting. We exploit the robustness of geometric methods to varying illumination to supervise our inverse rendering network. We apply a multiview stereo (MVS) pipeline to large sets of images of the same scene. We select pairs of overlapping images with different illumination, use the estimated relative pose and depth maps to cross project photometric invariants between views and use this for supervision via Siamese training. In other words, geometry provides correspondence that allows us to simulate varying illumination from a fixed viewpoint. Finally, the depth maps from MVS provide coarse normal map estimates that can be used for direct supervision of the normal map estimation.

1.1. Contribution

Deep learning has already shown good performance on components of the inverse rendering problem. This includes monocular depth estimation [11], depth and normal estimation [10] and intrinsic image decomposition [30]. However, these works use supervised learning. For tasks where ground truth does not exist, such approaches must either train on synthetic data (in which case generalisation to the real world is not guaranteed) or generate pseudo ground truth using an existing method (in which case the network is just learning to replicate the performance of the existing method). Inverse rendering of outdoor, complex scenes is itself an unsolved problem and so reliable ground truth is not available and supervised learning cannot be used. In this context, we make the following contributions. To the best of our knowledge, we are the first to exploit MVS supervision for learning inverse rendering. Second, we are the first to tackle the most general version of the problem, considering arbitrary outdoor scenes and learning from real data, as opposed to restricting to a single object class [46] or using synthetic training data [53]. Third, we introduce a statistical model of spherical harmonic lighting in natural scenes that we use as a prior. Finally, the resulting network is the first to inverse render all of shape, reflectance and lighting in the wild and we perform the first evaluation in this setting.

2. Related work

Classical approaches Classical methods estimate intrinsic properties by fitting photometric or geometric models. Most methods require multiple images. From multiview images, a structure-from-motion/multiview stereo pipeline enables recovery of dense mesh models [14, 24] though illumination effects are baked into the texture. From images with fixed viewpoint but varying illumination photometric stereo can be applied. Variants consider statistical BRDF models [3], the use of outdoor time-lapse images [29] and

spatially-varying BRDFs [18]. Attempts to combine geometric and photometric methods are limited. Haber *et al.* [19] assume known geometry (which can be provided by MVS) and inverse render reflectance and lighting from community photo collections. Kim *et al.* [26] represents the state-of-the-art and again uses an MVS initialisation for joint optimisation of geometry, illumination and albedo. Some methods consider a single image setting. Jeson *et al.* [22] introduce a local-adaptive reflectance smoothness constraint for intrinsic image decomposition on texture-free input images which are acquired with a texture separation algorithm. Barron *et al.* [4] present SIRFS, a classical optimisation-based approach that recovers all of shape, illumination and albedo using a sophisticated combination of generic priors.

Deep depth prediction Direct estimation of shape alone using deep neural networks has attracted a lot of attention. Eigen *et al.* [10, 11] were the first to apply deep learning in this context. Subsequently, performance gains were obtained using improved architectures [28], post-processing with classical CRF-based methods [36, 50, 51] and using ordinal relationships for objects within the scenes [8, 13, 34]. Zheng *et al.* [53] use synthetic images for training but improve generalisation using a synthetic-to-real transform GAN. However, all of this work requires supervision by ground truth depth. An alternative branch of methods explore using self-supervision from augmented data. For example, binocular stereo pairs can provide a supervisory signal through consistency of cross projected images [15, 16, 25]. Alternatively, video data can provide a similar source of supervision [48, 49, 54]. Some of other work built from specific ways were proposed recently. Tulsiani *et al.* [47] use multiview supervision in a ray tracing network. While all these methods take single image input, Ji *et al.* [23] tackle the MVS problem itself using deep learning.

Deep intrinsic image decomposition Intrinsic image decomposition is a partial step towards inverse rendering. It decomposes an image into reflectance (albedo) and shading but does not separate shading into shape and illumination. Even so, the lack of ground truth training data makes this a hard problem to solve with deep learning. Recent work either uses synthetic training data and supervised learning [7, 12, 20, 30, 39] or self-supervision/unsupervised learning. Very recently, Li *et al.* [33] used uncontrolled time-lapse images allowing them to combine an image reconstruction loss with reflectance consistency between frames. This work was further extended using photorealistic, synthetic training data [32]. Ma *et al.* [38] also trained on time-lapse sequences and introduced a new gradient constraint which encourage better explanations for sharp changes caused by shading or reflectance. Baslamisli *et al.* [5] applied a similar gradient constraint while they used supervised training.

Shelhamer *et al.* [43] propose a hybrid approach where a CNN estimates a depth map which is used to constrain a classical optimisation-based intrinsic image estimation.

Deep inverse rendering To date, this topic has not received much attention. One line of work simplifies the problem by restricting to a single object class, e.g. faces [46], meaning that a statistical face model can constrain the geometry and reflectance estimates. This enables entirely self-supervised training. Shu *et al.* [45] extend this idea with an adversarial loss. Sengupta *et al.* [42] on the other hand, initialise with supervised training on synthetic data, and fine-tune their network in an unsupervised fashion on real images. Aittala *et al.* [1] restrict geometry to almost planar objects and lighting to a flash in the viewing direction under which assumptions they can obtain impressive results. More general settings have been considered including natural illumination [31]. Kulkarni *et al.* [27] show how to learn latent variables that correspond to extrinsic parameters allowing image manipulation. The only prior work we are aware of that tackles the full inverse rendering problem requires direct supervision [21, 35, 37]. Hence, it is not applicable to scene-level inverse rendering, only objects, and relies on synthetic data for training, limiting the ability of the network to generalise to real images.

3. Preliminaries

We assume that a perspective camera observes a scene, such that the projection from 3D world coordinates, (u, v, w) , to 2D image coordinates, (x, y) , is given by:

$$\lambda \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \mathbf{P} \begin{bmatrix} u \\ v \\ w \\ 1 \end{bmatrix}, \quad \mathbf{P} = \mathbf{K} [\mathbf{R} \quad \mathbf{t}], \quad \mathbf{K} = \begin{bmatrix} f & 0 & c_x \\ 0 & f & c_y \\ 0 & 0 & 1 \end{bmatrix}, \quad (1)$$

where λ is an arbitrary scale factor, $\mathbf{R} \in SO(3)$ a rotation matrix, $\mathbf{t} \in \mathbb{R}^3$ a translation vector, f the focal length and (c_x, c_y) the principal point.

The inverse rendered shape estimate could be represented in a number of ways. For example, many previous methods estimate a viewer-centred depth map. However, local reflectance, and hence appearance, is determined by surface orientation, i.e. the local surface normal direction. So, to render a depth map for self-supervision, we would need to compute the surface normal. From a perspective depth map $w(x, y)$, the surface normal direction is:

$$\bar{\mathbf{n}} = \begin{bmatrix} -fw_x(x, y) \\ -fw_y(x, y) \\ (x - c_x)w_x(x, y) + (y - c_y)w_y(x, y) + w(x, y) \end{bmatrix} \quad (2)$$

from which the unit length normal is given by: $\mathbf{n} = \bar{\mathbf{n}} / \|\bar{\mathbf{n}}\|$. The derivatives of the depth map in the image

plane, $w_x(x, y)$ and $w_y(x, y)$, can be approximated by finite differences. However, (2) requires knowledge of the intrinsic camera parameters. This would severely restrict the applicability of our method. For this reason, we choose to estimate a surface normal map directly.

Although the surface normal can be represented by a 3D vector, since $\|\mathbf{n}\| = 1$ it has only two degrees of freedom. So, our network estimates the two elements of the surface gradient at each pixel, $w_u(x, y)$ and $w_v(x, y)$, and the transformation to a 3D surface normal vector is computed by a fixed layer that calculates: $\bar{\mathbf{n}} = [-w_u(x, y), -w_v(x, y), 1]^T$. Note that we estimate the normal map in a viewer-centred coordinate system.

We assume that appearance can be approximated by a local reflectance model under environment illumination. Specifically we use a Lambertian diffuse model with order 2 spherical harmonic lighting. This means that RGB intensity can be computed as

$$\mathbf{i}_{\text{lin}}(\mathbf{n}, \boldsymbol{\alpha}, \mathbf{L}) = \text{diag}(\boldsymbol{\alpha})\mathbf{L}\mathbf{b}(\mathbf{n}), \quad (3)$$

where $\mathbf{L} \in \mathbb{R}^{3 \times 9}$ contains the spherical harmonic colour illumination coefficients, $\boldsymbol{\alpha} = [\alpha_r, \alpha_g, \alpha_b]^T$ is the colour diffuse albedo and the order 2 basis is given by:

$$\mathbf{b}(\mathbf{n}) = [1, n_x, n_y, n_z, 3n_z^2 - 1, n_x n_y, n_x n_z, n_y n_z, n_x^2 - n_y^2]^T. \quad (4)$$

Our appearance model means that we neglect high frequency illumination effects, cast shadows and interreflections. However, we found that in practice this model works well for typical outdoor scenes. Finally, cameras apply a nonlinear gamma transformation. We simulate this to produce our final predicted intensities: $\mathbf{i}_{\text{pred}} = \mathbf{i}_{\text{lin}}^{1/\gamma}$, where we assume a fixed $\gamma = 2.2$.

4. Architecture

Our inverse rendering network (see Fig. 2) is an image-to-image network that regresses albedo and normal maps from a single image and uses these to estimate lighting. We describe these inference components in more detail here.

4.1. Trainable encoder-decoder

We implement a deep fully-convolutional neural network with skip connections like the hourglass architecture [41]. We use a single encoder and separate deconvolution decoders for albedo and normal prediction. Albedo maps have 3 channel RGB output, normal maps have two channels for the surface gradient which is converted to a normal map as described above. Both convolutional subnet and deconvolutional subnet contain 15 layers and the activation functions are ReLUs. Adam Optimiser is used in training.

4.2. Implicit lighting prediction

In order to estimate illumination parameters, one option would be to use a fully connected branch from the output

of our decoder and train our network to predict it directly. However, fully connected layers require very large numbers of parameters and, in fact, lighting can be inferred from the input image and estimated albedo and normal maps, making its explicit prediction redundant. An additional advantage is that the architecture remains fully convolutional and so can process images of any size at inference time.

Consider an input image comprising K pixels. We invert the nonlinear gamma and stack the linearised RGB values to form the matrix $\mathbf{I} \in \mathbb{R}^{3 \times K}$. We similarly stack the estimated albedo map to form $\mathbf{A} \in \mathbb{R}^{3 \times K}$, the estimated surface normals to form $\mathbf{N} \in \mathbb{R}^{3 \times K}$ and define $\mathbf{B}(\mathbf{N}) \in \mathbb{R}^{9 \times K}$ by applying (4) to each normal vector. We can now rewrite (3) for the whole image as:

$$\mathbf{I} = \mathbf{A} \odot \mathbf{L}\mathbf{B}(\mathbf{N}), \quad (5)$$

where \odot is the Hadamard (element-wise) product. We can now solve for the spherical harmonic illumination coefficients in a least squares sense, using the whole image. This can be done using any method, so long as the computation is differentiable such that losses dependent on the estimated illumination can have their gradients backpropagated into the inverse rendering network. For example, the solution using the pseudoinverse is given by: $\mathbf{L} = (\mathbf{I} \oslash \mathbf{A})\mathbf{B}(\mathbf{N})^+$, where \oslash denotes element-wise division and $\mathbf{B}(\mathbf{N})^+$ is the pseudoinverse of $\mathbf{B}(\mathbf{N})$. Fig. 2 shows the inferred shading, $\mathbf{I} \oslash \mathbf{A}$, and a visualisation of the estimated lighting.

5. Supervision

As shown in Fig. 2, we use a data term (the error between predicted and observed appearance) for self-supervision. However, inverse rendering using only a data term is ill-posed (an infinite set of solutions can yield zero data error) and so we use additional sources of supervision, all of which are essential for good performance. We describe all sources of supervision in this section.

5.1. Self-supervision via differentiable rendering

Given estimated normal and albedo maps and spherical harmonic illumination coefficients, we compute a predicted image using (3). This local illumination model is straightforward to differentiate. Self-supervision is provided by the error between the predicted, \mathbf{i}_{pred} , and observed, \mathbf{i}_{obs} , intensities. We compute this error in LAB space as this provides perceptually more convincing results:

$$\ell_{\text{appearance}} = \|\text{LAB}(\mathbf{i}_{\text{pred}}) - \text{LAB}(\mathbf{i}_{\text{obs}})\|, \quad (6)$$

where LAB performs the colour space transformation.

5.2. Natural illumination model and prior

The spherical harmonic lighting model in (3) enables efficient representation of complex lighting. However, even

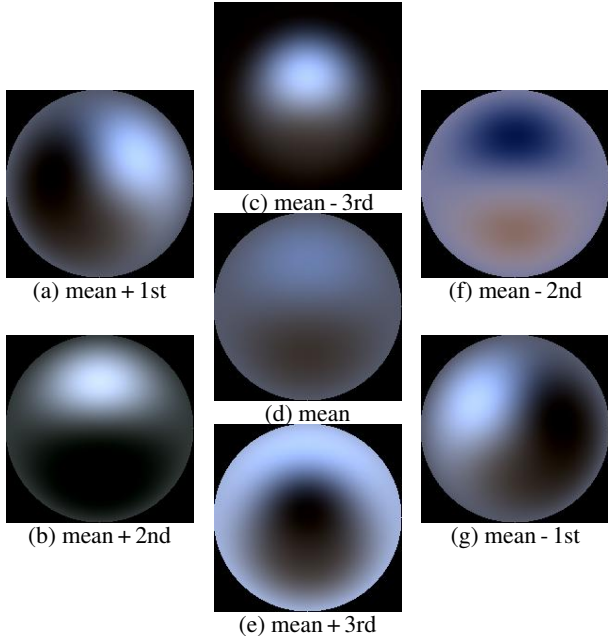


Figure 3: Statistical illumination model. The central image shows the mean illumination. The two diagonals and the vertical show the first 3 principal components.

within this low dimensional space, not all possible illumination environments are natural. The space of natural illumination possesses statistical regularities [9]. We can use this knowledge to constrain the space of possible illumination and enforce a prior on the illumination parameters. To do this, we build a statistical illumination model (see Fig. 3) using a dataset of 79 HDR spherical panoramic images taken outdoors. For each environment, we compute the spherical harmonic coefficients, $\mathbf{L}_i \in \mathbb{R}^{3 \times 9}$. Since the overall intensity scale is arbitrary, we also normalise each lighting matrix to unit norm, $\|\mathbf{L}_i\|_{\text{Fro}} = 1$, to avoid ambiguity with the albedo scale. Our illumination model in (5) uses surface normals in a viewer-centred coordinate system. So, the dataset must be augmented to account for possible rotations of the environment relative to the viewer. Since the rotation around the vertical (v) axis is arbitrary, we rotate the lighting coefficients by angles between 0 and 2π in increments of $\pi/18$. In addition, to account for camera pitch or roll, we additionally augment with rotations about the u and w axes in the range $(-\pi/6, \pi/6)$. This gives us a dataset of 139,356 environments. We then build a statistical model, such that any illumination can be approximated as:

$$\text{vec}(\mathbf{L}) = \mathbf{P} \text{diag}(\sigma_1, \dots, \sigma_D) \boldsymbol{\alpha} + \text{vec}(\bar{\mathbf{L}}). \quad (7)$$

where $\mathbf{P} \in \mathbb{R}^{27 \times D}$ contains the principal components, $\sigma_1^2, \dots, \sigma_D^2$ are the corresponding eigenvalues, $\bar{\mathbf{L}} \in \mathbb{R}^{3 \times 9}$ is the mean lighting coefficients and $\boldsymbol{\alpha} \in \mathbb{R}^D$ is the parametric representation of \mathbf{L} . We use $D = 18$ dimensions. Under the assumption that the original data is Gaussian distributed

then the parameters are normally distributed: $\boldsymbol{\alpha} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. When we compute lighting, we do so within the subspace of the statistical model. In addition, we introduce a prior loss on the estimated lighting vector: $\ell_{\text{lighting}} = \|\boldsymbol{\alpha}\|^2$.

5.3. Multiview stereo supervision

A pipeline comprising structure-from-motion followed by multiview stereo (which we refer to simply as MVS) enables both camera poses and dense 3D scene models to be estimated from large, uncontrolled image sets. Of particular importance to us, these pipelines are relatively insensitive to illumination variation between images in the dataset since they rely on matching local image features that are themselves illumination insensitive. We emphasise that MVS is run offline prior to training and that at inference time our network uses only single images of novel scenes. We use the MVS output for three sources of supervision.

Cross-projection We use the MVS poses and depth maps to establish correspondence between views, allowing us to cross-project quantities between overlapping images. Given an estimated depth map, $w(x, y)$, in view i and camera matrices for views i and j , a pixel (x, y) can be cross-projected to location (x', y') in view j via:

$$\lambda \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = \mathbf{P}_j \begin{bmatrix} \mathbf{R}_i^T & -\mathbf{R}_i^T \mathbf{t}_i \\ \mathbf{0} & 1 \end{bmatrix} \begin{bmatrix} w(x, y) \mathbf{K}_i^{-1} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \\ 1 \end{bmatrix} \quad (8)$$

In practice, we perform the cross-projection in the reverse direction, computing non-integer pixel locations in the source view for each pixel in the target view. We can then use bilinear interpolation of the source image to compute quantities for each pixel in the target image. Since the MVS depth maps contain holes, any pixels that cross project to a missing pixel are not assigned a value. Similarly, any target pixels that project outside the image bounds of the source are not assigned a value.

Direct normal map supervision The per-view depth maps provided by MVS can be used to estimate normal maps, albeit that they are typically coarse and incomplete (see Fig. 1, column 5). We compute guide normal maps from the depth maps and intrinsic camera parameters estimated by MVS using (2). The guide normal maps are used for direct supervision by computing a loss that measures the angular difference between the guide, $\mathbf{n}_{\text{guide}}$, and estimated, \mathbf{n}_{est} , surface normals: $\ell_{\text{NM}} = \arccos(\mathbf{n}_{\text{guide}} \cdot \mathbf{n}_{\text{est}})$.

Albedo consistency loss Diffuse albedo is an intrinsic quantity. Hence, we expect that albedo estimates of the same scene point from two overlapping images should be the same, even if the illumination varies between views. Hence, we automatically select pairs of images that overlap (defined as having similar camera locations and similar

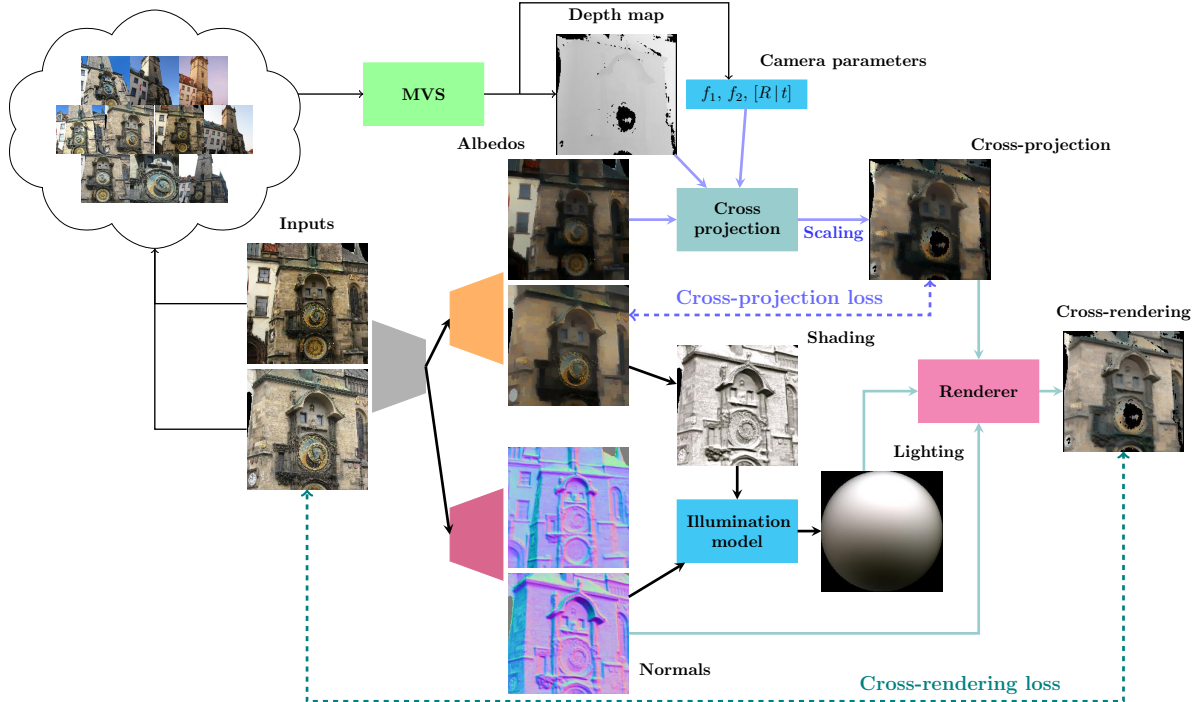


Figure 4: Siamese MVS supervision: albedo cross-projection consistency and cross-rendering losses (shown in one direction for simplicity). Note: shading depends on input and albedo as in Fig. 2 but this dependency is excluded for simplicity.

centres of mass of their backprojected depth maps). We discard pairs that do not contain illumination variation (where cross-projected appearance is too similar). Then, we train our network in a Siamese fashion on these pairs and use the cross projection described above to compute an albedo consistency loss: $\ell_{\text{albedo}} = \|\text{LAB}(\mathbf{A}_i) - \text{LAB}(s\mathbf{A}_j)\|_{\text{fro}}^2$, where $\mathbf{A}_i, \mathbf{A}_j \in \mathbb{R}^{3 \times K}$ are the estimated albedo maps in the i th and j th images respectively, where \mathbf{A}_j has been cross projected to view i , for the K pixels in which image i has a defined depth value. The scalar s is the value that minimises the loss and accounts for the fact that there is an overall scale ambiguity between images. Again, we compute albedo consistency loss in LAB space. The albedo consistency loss is visualised by the blue arrows in Fig. 4.

Cross-rendering loss For improved stability, we also use a mixed cross-projection/appearance loss, $\ell_{\text{cross-rend}}$. We use the cross-projected albedo above in conjunction with the estimated normals and illumination to render a new image and measure the appearance error in the same way as (6). This loss is visualised by the green arrows in Fig. 4.

5.4. Albedo priors

Finally, we also employ two additional prior losses on the albedo. This helps resolve ambiguities between shading and albedo. First, we introduce an albedo smoothness prior, $\ell_{\text{albedo-smooth}}$. Rather than uniformly applying smoothness

penalty, we apply a pixel-wise varying weighted penalty according to chromaticities of the input image. So the stronger smoothness penalties are only enforced on neighbouring pixels with closer chromaticities. The loss itself is the L1 distance between adjacent pixels.

Second, during the self-supervised phase of training, we also introduce a pseudo supervision loss to prevent convergence to trivial solutions. After the pretraining process (see Section 6), our model learns plausible albedo predictions using MVS normals. To prevent subsequent training diverging too far from this, we encourage albedo predictions to remain close to the pretrained albedo predictions.

6. Training

We train our network to minimise: $\ell = w_1 \ell_{\text{appearance}} + w_2 \ell_{\text{NM}} + w_3 \ell_{\text{albedo}} + w_4 \ell_{\text{cross-rend}} + w_5 \ell_{\text{albedo-smooth}} + w_6 \ell_{\text{albedo-pseudoSup}}$.

Datasets We train using the MegaDepth [34] dataset. This contains dense depth maps and camera calibration parameters estimated from crawled Flickr images. The pre-processed images have arbitrary shapes and orientations. For ease of training, we crop square images and resize to a fixed size. We choose our crops to maximise the number of pixels with defined depth values. Where possible, we crop multiple images from each image, achieving augmen-

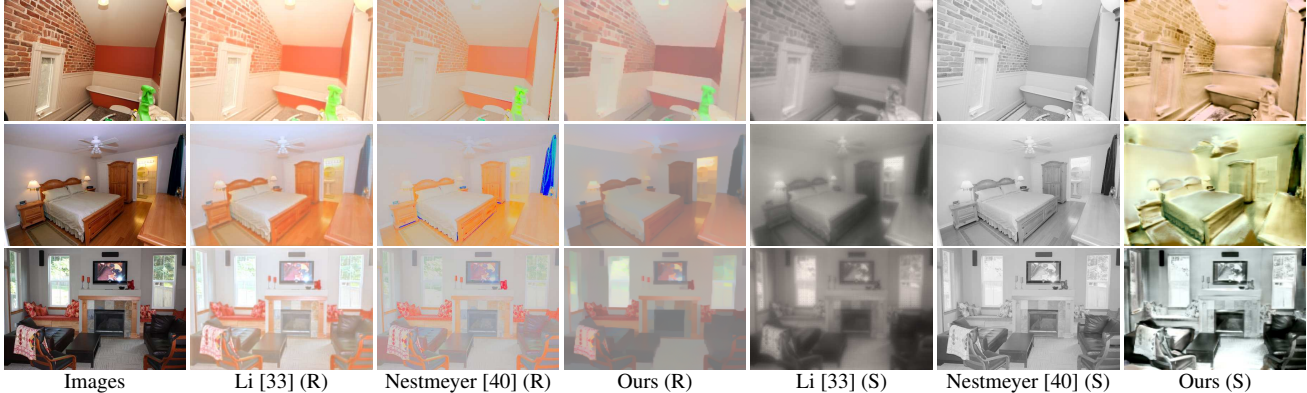


Figure 5: Qualitative results for I IW. Second column to forth column are reflectance predictions from [33], [40] and ours. The last three columns are corresponding shading predictions.

tation as well as standardisation. We create mini-batches with overlap between all pairs of images in the mini-batch and sufficient illumination variation (correlation coefficient of intensity histograms significantly different from 1). Finally, before inputting an image to our network, we detect and mask the sky region using PSPNet [52]. This is because the albedo map and normal map in sky area are meaningless and it severely influences illumination estimation.

Training strategy We found that for convergence to a good solution it is important to include a pre-training phase. During this phase, the surface normals used for illumination estimation and for the appearance-based losses are the MVS normal maps. This means that the surface normal prediction decoder is only learning from the direct supervision loss, i.e. it is learning to replicate the MVS normals. After this initial phase, we switch to full self-supervision where the predicted appearance is computed entirely from estimated quantities. Note that this pre-training step is not using pseudo albedo supervisions.

7. Evaluation

There are no existing benchmarks for inverse rendering in the wild. So, we evaluate our method on an intrinsic image benchmark and devise our own benchmark for inverse rendering. Finally, we show a relighting application.

Evaluation on I IW The standard benchmark for intrinsic image decomposition is Intrinsic Images in the Wild [6] (I IW) which is almost exclusively indoor scenes. Since our training regime requires large multiview image datasets, we are restricted to using scene-tagged images crawled from the web, which are usually outdoors. In addition, our illumination model is learnt on outdoor, natural environments. For these reasons, we cannot perform training or fine-tuning on indoor benchmarks. Moreover, our network is not trained specifically for the task of intrinsic image estimation and

Methods	Training data	WHDR
Nestmeyer [40] (CNN)	I IW	19.5
Zhou <i>et al.</i> [55]	I IW	19.9
Fan <i>et al.</i> [12]	I IW	14.5
DI [39]	Sintel+MIT	37.3
Shi <i>et al.</i> [44]	ShapeNet	59.4
Li <i>et al.</i> [33]	BigTime	20.3
Ours	MegaDepth	21.4

Table 1: Evaluation results on I IW benchmark using WHDR percentage (lower is better). The second column shows which dataset on which the networks were trained.

our shading predictions are limited by the fact that we use an explicit local illumination model (so cannot predict cast shadows). Nevertheless, we test our network on this benchmark directly without fine-tuning. We follow the suggestion in [40] and rescale albedo predictions to the range (0.5, 1) before evaluation. Quantitative results are shown in Tab. 1 and some qualitative visual comparison in Fig. 5. Despite the limitations described above, we achieve the second best performance of the methods not trained on the I IW data.

Evaluation on MegaDepth We evaluate inverse rendering using unobserved scenes from the MegaDepth dataset [34]. We evaluate normal estimation performance directly using the MVS geometry. We evaluate albedo estimation using a state-of-the-art multiview inverse rendering algorithm [26]. Given the output from their pipeline, we perform rasterisation to generate albedo ground truth for every input image. Note that both sources of “ground truth” here are themselves only estimations, e.g. the albedo ground truth contains ambient occlusion baked in. The colour balance of the estimated albedo is arbitrary, so we compute per-channel optimal scalings prior to computing errors. We use three metrics - MSE, LMSE and DSSIM, which are commonly used for evaluating albedo predictions. To eval-

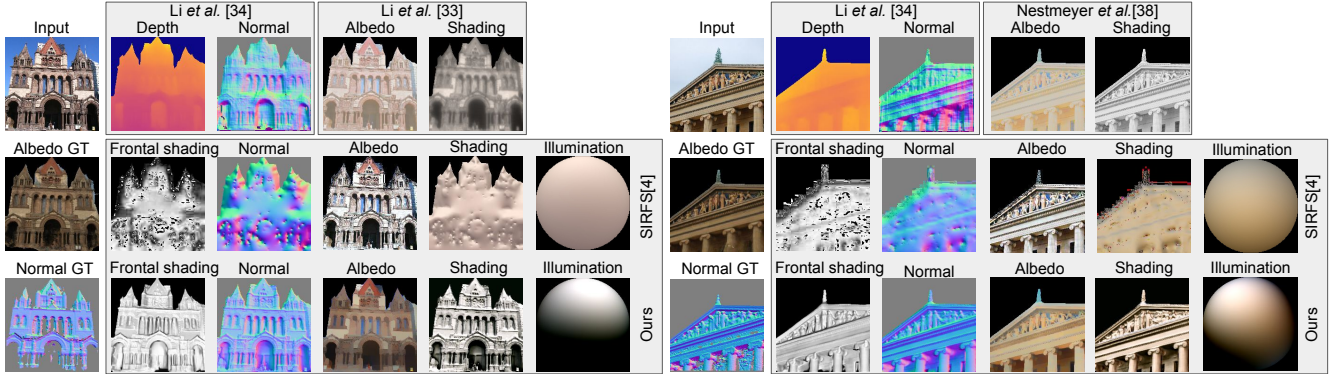


Figure 6: Inverse Rendering Results.

Methods	Reflectances			Normals	
	MSE	LMSE	DSSIM	Mean	Median
Li <i>et al.</i> [34]	-	-	-	50.6	50.4
Godard <i>et al.</i> [16]	-	-	-	79.2	79.6
Nestmeyer <i>et al.</i> [40]	0.0204	0.0735	0.241	-	-
Li <i>et al.</i> [33]	0.0171	0.0637	0.208	-	-
SIRFS [4]	0.0383	0.222	0.270	50.6	48.5
Ours	0.0170	0.0718	0.201	37.7	34.8

Table 2: Quantitative inverse rendering results. Reflectance (albedo) errors are measured against multiview inverse rendering result [26] and normals against MVS results.

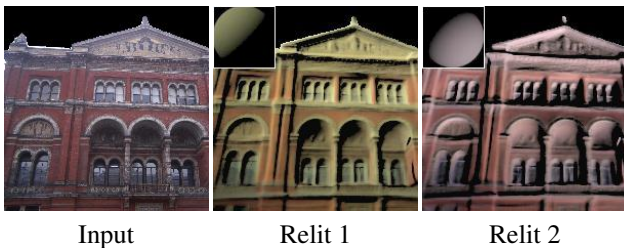


Figure 7: Relighting results from predicted albedo and normal maps (see Fig. 1, row 3). The novel lighting is shown in the upper left corner.

uate normal predictions, we use angular errors. The correctness of illumination predictions could be inferred by the other two, so we do not perform explicit evaluations on it. The quantitative evaluations are shown in Tab. 2. For depth prediction methods, we first compute the optimal scaling onto the ground truth geometry, then differentiate to compute surface normals. These methods can only be evaluated on normal prediction. Intrinsic image methods can only be evaluated on albedo prediction. We can see that our network performs best in normal prediction and also the best in MSE and DSSIM. Qualitative example results can be seen in Fig. 6.

Relighting Finally, as an example application we show

that our inverse rendering result is sufficiently stable for realistic relighting. A scene from Fig. 1 is relit in Fig. 7 with two novel illuminations. Both show realistic shading and overall colour balance.

8. Conclusions

We have shown for the first time that the task of inverse rendering can be learnt from real world images in uncontrolled conditions. Our results show that “shape-from-shading” in the wild is possible and are far superior to classical methods. It is interesting to ponder how this feat is achieved. We believe the reason this is possible is because of the large range of cues that the deep network can exploit, for example shading, texture, ambient occlusion, perhaps even high level semantic concepts learnt from the diverse data. For example, once a region is recognised as a “window”, the possible shape and configuration is much restricted. Recognising a scene as a man-made building suggests the presence of many parallel and orthogonal planes. These sort of cues would be extremely difficult to exploit in hand-crafted solutions.

There are many promising ways in which this work can be extended. First, our modelling assumptions could be relaxed, for example using more general reflectance models and estimating global illumination effects such as shadowing. Second, our network could be combined with a depth prediction network. Either the two networks could be applied independently and then the depth and normal maps merged, or a unified network could be trained in which the normals computed from the depth map are used to compute the losses we use in this paper. Third, our network could benefit from losses used in training intrinsic image decomposition networks. For example, if we added the timelapse dataset of [33] to our training, we could incorporate their reflectance consistency loss to improve our albedo map estimates. Our code, trained model and inverse rendering benchmark data is available at <https://github.com/YeeU/InverseRenderNet>.

References

- [1] Miika Aittala, Timo Aila, and Jaakko Lehtinen. Reflectance modeling by neural texture synthesis. *ACM Transactions on Graphics (TOG)*, 35(4):65, 2016.
- [2] O Aldrian and WA Smith. Inverse rendering of faces with a 3d morphable model. *IEEE transactions on pattern analysis and machine intelligence*, 35(5):1080–1093, 2013.
- [3] N. Alldrin, T. Zickler, and D. Kriegman. Photometric stereo with non-parametric and spatially-varying reflectance. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, June 2008.
- [4] Jonathan T Barron and Jitendra Malik. Shape, illumination, and reflectance from shading. *TPAMI*, 2015.
- [5] Anil S. Baslamisli, Hoang-An Le, and Theo Gevers. Cnn based learning using reflection and retinex models for intrinsic image decomposition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [6] Sean Bell, Kavita Bala, and Noah Snavely. Intrinsic images in the wild. *ACM Trans. on Graphics (SIGGRAPH)*, 33(4), 2014.
- [7] Sai Bi, Nima Khademi Kalantari, and Ravi Ramamoorthi. Deep Hybrid Real and Synthetic Training for Intrinsic Decomposition. In Wenzel Jakob and Toshiya Hachisuka, editors, *Eurographics Symposium on Rendering - Experimental Ideas & Implementations*. The Eurographics Association, 2018.
- [8] Weifeng Chen, Zhao Fu, Dawei Yang, and Jia Deng. Single-image depth perception in the wild. In *Advances in Neural Information Processing Systems*, pages 730–738, 2016.
- [9] Ron O Dror, Thomas K Leung, Edward H Adelson, and Alan S Willsky. Statistics of real-world illumination. In *Proc. CVPR*, 2001.
- [10] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2650–2658, 2015.
- [11] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Advances in neural information processing systems*, pages 2366–2374, 2014.
- [12] Qingnan Fan, Jiaolong Yang, Gang Hua, Baoquan Chen, and David Wipf. Revisiting deep intrinsic image decompositions. In *Proceedings of The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8944–8952, 2018.
- [13] Huan Fu, Mingming Gong, Chaohui Wang, Kayhan Batmanghelich, and Dacheng Tao. Deep ordinal regression network for monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2002–2011, 2018.
- [14] Yasutaka Furukawa and Jean Ponce. Accurate, dense, and robust multiview stereopsis. *IEEE Trans. Pattern Anal. Mach. Intell.*, 32(8):1362–1376, Aug. 2010.
- [15] Ravi Garg, Vijay Kumar BG, Gustavo Carneiro, and Ian Reid. Unsupervised cnn for single view depth estimation: Geometry to the rescue. In *European Conference on Computer Vision*, pages 740–756. Springer, 2016.
- [16] Clément Godard, Oisín Mac Aodha, and Gabriel J Brostow. Unsupervised monocular depth estimation with left-right consistency. In *CVPR*, volume 2, page 7, 2017.
- [17] Michael Goesele, Noah Snavely, Brian Curless, Hugues Hoppe, and Steven M. Seitz. Multi-view stereo for community photo collections. *2007 IEEE 11th International Conference on Computer Vision*, pages 1–8, 2007.
- [18] Dan B Goldman, Brian Curless, Aaron Hertzmann, and Steven M Seitz. Shape and spatially-varying brdfs from photometric stereo. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 32(6):1060–1071, 2010.
- [19] T. Haber, C. Fuchs, P. Bekaer, H. P. Seidel, M. Goesele, and H. P. A. Lensch. Relighting objects from image collections. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 627–634, June 2009.
- [20] Guangyun Han, Xiaohua Xie, Jianhuang Lai, and Wei-Shi Zheng. Learning an intrinsic image decomposer using synthesized rgb-d dataset. *IEEE Signal Processing Letters*, 25(6):753–757, 2018.
- [21] Michael Janner, Jiajun Wu, Tejas D Kulkarni, Ilker Yildirim, and Josh Tenenbaum. Self-supervised intrinsic image decomposition. In *Advances in Neural Information Processing Systems*, pages 5936–5946, 2017.
- [22] Junho Jeon, Sunghyun Cho, Xin Tong, and Seungyong Lee. Intrinsic image decomposition using structure-texture separation and surface normals. In *European Conference on Computer Vision*, pages 218–233. Springer, 2014.
- [23] Mengqi Ji, Juergen Gall, Haitian Zheng, Yebin Liu, and Lu Fang. SurfacerNet: an end-to-end 3d neural network for multiview stereopsis. *arXiv preprint arXiv:1708.01749*, 2017.
- [24] Michael Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM Trans. Graph.*, 32(3):29:1–29:13, July 2013.
- [25] Alex Kendall, Hayk Martirosyan, Saumitro Dasgupta, Peter Henry, Ryan Kennedy, Abraham Bachrach, and Adam Bry. End-to-end learning of geometry and context for deep stereo regression. *CoRR*, vol. abs/1703.04309, 2017.
- [26] Kichang Kim, Akihiko Torii, and Masatoshi Okutomi. Multi-view inverse rendering under arbitrary illumination and albedo. In *European Conference on Computer Vision*, pages 750–767. Springer, 2016.
- [27] Tejas D Kulkarni, William F Whitney, Pushmeet Kohli, and Josh Tenenbaum. Deep convolutional inverse graphics network. In *Advances in neural information processing systems*, pages 2539–2547, 2015.
- [28] Iro Laina, Christian Rupprecht, Vasileios Belagiannis, Federico Tombari, and Nassir Navab. Deeper depth prediction with fully convolutional residual networks. In *3D Vision (3DV), 2016 Fourth International Conference on*, pages 239–248. IEEE, 2016.
- [29] Fabian Langguth. Photometric stereo for outdoor webcams. In *Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, CVPR '12, pages 262–269, Washington, DC, USA, 2012. IEEE Computer Society.
- [30] Louis Lettry, Kenneth Vanhoey, and Luc Van Gool. Darn: a deep adversarial residual network for intrinsic image decom-

- position. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1359–1367. IEEE, 2018.
- [31] Xiao Li, Yue Dong, Pieter Peers, and Xin Tong. Modeling surface appearance from a single photograph using self-augmented convolutional neural networks. *ACM Transactions on Graphics (TOG)*, 36(4):45, 2017.
- [32] Zhengqi Li and Noah Snavely. Cgintrinsics: Better intrinsic image decomposition through physically-based rendering. In *European Conference on Computer Vision (ECCV)*, 2018.
- [33] Zhengqi Li and Noah Snavely. Learning intrinsic image decomposition from watching the world. In *Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [34] Zhengqi Li and Noah Snavely. Megadepth: Learning single-view depth prediction from internet photos. In *Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [35] Zhengqi Li, Zexiang Xu, Ravi Ramamoorthi, Kalyan Sunkavalli, and Manmohan Chandraker. Learning to reconstruct shape and spatially-varying reflectance from a single image. In *SIGGRAPH Asia 2018 Technical Papers*, page 269. ACM, 2018.
- [36] Fayao Liu, Chunhua Shen, and Guosheng Lin. Deep convolutional neural fields for depth estimation from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5162–5170, 2015.
- [37] Guilin Liu, Duygu Ceylan, Ersin Yumer, Jimei Yang, and Jyh-Ming Lien. Material editing using a physically based rendering network. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2261–2269, 2017.
- [38] Wei-Chiu Ma, Hang Chu, Bolei Zhou, Raquel Urtasun, and Antonio Torralba. Single image intrinsic decomposition without a single intrinsic image. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 201–217, 2018.
- [39] Takuya Narihira, Michael Maire, and Stella X Yu. Direct intrinsics: Learning albedo-shading decomposition by convolutional regression. In *Proceedings of the IEEE international conference on computer vision*, pages 2992–2992, 2015.
- [40] Thomas Nestmeyer and Peter V Gehler. Reflectance adaptive filtering improves intrinsic image estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 2, page 4, 2017.
- [41] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *Proc. ECCV*, pages 483–499, 2016.
- [42] Soumyadip Sengupta, Angjoo Kanazawa, Carlos D Castillo, and David W Jacobs. Sfsnet: Learning shape, reflectance and illuminance of faces in the wild. *arXiv preprint arXiv:1712.01261*, 2017.
- [43] Evan Shelhamer, Jonathan T Barron, and Trevor Darrell. Scene intrinsics and depth from a single image. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 37–44, 2015.
- [44] Jian Shi, Yue Dong, Hao Su, and X Yu Stella. Learning non-lambertian object intrinsics across shapenet categories. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 5844–5853. IEEE, 2017.
- [45] Zhixin Shu, Ersin Yumer, Sunil Hadap, Kalyan Sunkavalli, Eli Shechtman, and Dimitris Samaras. Neural face editing with intrinsic image disentangling. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 5444–5453. IEEE, 2017.
- [46] Ayush Tewari, Michael Zollhöfer, Hyeonwoo Kim, Pablo Garrido, Florian Bernard, Patrick Pérez, and Christian Theobalt. Mofa: Model-based deep convolutional face autoencoder for unsupervised monocular reconstruction. In *The IEEE International Conference on Computer Vision (ICCV)*, volume 2, page 5, 2017.
- [47] Shubham Tulsiani, Tinghui Zhou, Alexei A Efros, and Jitendra Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. In *CVPR*, volume 1, page 3, 2017.
- [48] Sudheendra Vijayanarasimhan, Susanna Ricco, Cordelia Schmid, Rahul Sukthankar, and Katerina Fragkiadaki. Sfmmnet: Learning of structure and motion from video. *arXiv preprint arXiv:1704.07804*, 2017.
- [49] Chaoyang Wang, José Miguel Buenaposada, Rui Zhu, and Simon Lucey. Learning depth from monocular videos using direct methods. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2022–2030, 2018.
- [50] Peng Wang, Xiaohui Shen, Zhe Lin, Scott Cohen, Brian Price, and Alan L Yuille. Towards unified depth and semantic prediction from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2800–2809, 2015.
- [51] Dan Xu, Elisa Ricci, Wanli Ouyang, Xiaogang Wang, and Nicu Sebe. Multi-scale continuous crfs as sequential deep networks for monocular depth estimation. In *Proceedings of CVPR*, volume 1, 2017.
- [52] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 2881–2890, 2017.
- [53] Chuanxia Zheng, Tat-Jen Cham, and Jianfei Cai. T2net: Synthetic-to-realistic translation for solving single-image depth estimation tasks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 767–783, 2018.
- [54] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *CVPR*, volume 2, page 7, 2017.
- [55] Tinghui Zhou, Philipp Krahenbuhl, and Alexei A Efros. Learning data-driven reflectance priors for intrinsic image decomposition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3469–3477, 2015.