# Self-calibrating Deep Photometric Stereo Networks

Guanying Chen[1]    Kai Han[2]    Boxin Shi[3,4]    Yasuyuki Matsushita[5]    Kwan-Yee K. Wong[1]

[1]The University of Hong Kong    [2]University of Oxford
[3]Peking University    [4]Peng Cheng Laboratory    [5]Osaka University

## Abstract

*This paper proposes an uncalibrated photometric stereo method for non-Lambertian scenes based on deep learning. Unlike previous approaches that heavily rely on assumptions of specific reflectances and light source distributions, our method is able to determine both shape and light directions of a scene with unknown arbitrary reflectances observed under unknown varying light directions. To achieve this goal, we propose a two-stage deep learning architecture, called* SDPS-Net*, which can effectively take advantage of intermediate supervision, resulting in reduced learning difficulty compared to a single-stage model. Experiments on both synthetic and real datasets show that our proposed approach significantly outperforms previous uncalibrated photometric stereo methods.*

## 1. Introduction

Photometric stereo aims at recovering the surface normal of a static object from a set of images captured under different light directions [34, 29]. *Calibrated* photometric stereo methods assume known light directions, and promising results have been reported [28] at the cost of tedious light source calibration. The problem of *uncalibrated* photometric stereo, where light directions are unknown, still remains an open challenge, and its stable solution is wanted because of the ease of setting. In this work, we study the problem of uncalibrated photometric stereo for surfaces with general and unknown isotropic reflectance.

Most of the existing methods for uncalibrated photometric stereo [2, 27, 23] assume a simplified reflectance model, such as the Lambertian model, and focus on resolving the shape-light ambiguity, such as the Generalized Bas-Relief (GBR) ambiguity [3]. Although methods of [19, 20] can handle surfaces with general bidirectional reflectance distribution functions (BRDFs), they rely on a uniform distribution of light directions for deriving a solution.

Recently, with the great success of deep learning in various computer vision tasks, deep learning based methods have been introduced to calibrated photometric stereo [25, 31, 15, 5]. Instead of explicitly modeling complex surface reflectances, they directly learn the mapping from reflectance observations to surface normals given light directions. Although they have obtained promising results in a calibrated setting, they cannot handle the more challenging problem of *uncalibrated* photometric stereo, where light directions are unknown. One simple strategy to handle uncalibrated photometric stereo with deep learning is to directly learn the mapping from images to surface normals without taking the light directions as input. However, as reported in [5], the performance of such a model lags far behind those which take both images and light directions as input.

In this paper, we propose a two-stage model named Self-calibrating Deep Photometric Stereo Networks (SDPS-Net) to tackle this problem. The first stage of SDPS-Net, denoted as *Lighting Calibration Network* (LCNet), takes an arbitrary number of images as input and estimates their corresponding light directions and intensities. The second stage of SDPS-Net, denoted as *Normal Estimation Network* (NENet), estimates a surface normal map of a scene based on the lighting conditions estimated by LCNet and the input images. The rationales behind the design of our two-stage model are as follows. First, lighting information is very important for normal estimation since lighting is the source of various cues, such as shading and reflectance, and estimating the light directions (3-vectors) and intensities (scalars) is in principle much easier than directly estimating the normal map (a 3-vector at each pixel location) together with the lighting conditions. Second, by explicitly learning to estimate light directions and intensities, the model can take advantage of the intermediate supervision, resulting in a more interpretable behavior. Last, the proposed LCNet can be seamlessly integrated with existing calibrated photometric stereo methods, which enables them to deal with unknown lighting conditions. Our code and model can be found at https://guanyingc.github.io/SDPS-Net.

## 2. Related Work

In this section, we review learning based photometric stereo and uncalibrated photometric stereo methods. We also briefly review the loosely related work on learning

based lighting estimation. Readers are referred to [28] for a comprehensive survey on calibrated photometric stereo with Lambertian surfaces and general BRDFs using non-learning based methods.

**Learning based photometric stereo** Recently, a few deep learning based methods have been introduced to calibrated photometric stereo [25, 31, 15, 5]. Santo *et al.* [25] proposed a fully-connected network to learn the mapping from reflectance observations captured under a pre-defined set of light directions to surface normal in a pixel-wise manner. Taniai and Maehara [31] introduced an unsupervised learning framework that predicts both the surface normals and reflectance images of an object. Their model is "trained" at test time for each test object by minimizing the reconstruction loss between the input images and the rendered images. Ikehata [15] introduced a fixed shape representation, called observation map, that is invariant to the number and permutation of the images. For each surface point of the object, all its observations are merged into an observation map based on the given light directions, and the observation map is then fed to a convolutional neural network (CNN) to regress the normal vector. Chen *et al.* [5] proposed a fully-convolutional network (FCN) to infer the normal map from the input image-lighting pairs, and an order-agnostic max-pooling operation was adopted to handle an arbitrary number of inputs. All the above methods assume known lighting conditions and cannot handle uncalibrated photometric stereo, where the light directions and intensities are not known a priori.

**Uncalibrated photometric stereo** When lighting is unknown, the surface normals of a Lambertian object can only be estimated up to a $3 \times 3$ linear ambiguity [12], which can be reduced to a 3-parameter GBR ambiguity [3, 36] using the surface integrability constraint. Previous work used additional clues like albedo priors [2, 27], inter-reflections [4], specular spikes [7], Torrance and Sparrow reflectance model [11], reflectance symmetry [30, 35], multi-view images [9], and local diffuse maxima [23], to resolve the GBR ambiguity. Cho *et al.* [6] considered a semi-calibrated case where the light directions are known but not their intensities. There are few works that can handle non-Lambertian surfaces under unknown lighting. Hertzmann and Seitz [13] proposed an exemplar based method by inserting an additional reference object to the scene. Methods based on cues like similarity in radiance changes [26, 19] and attached shadow [22] were also introduced, but they require the light sources to be uniformly distributed on the whole sphere. Recently, Lu *et al.* [18] introduced a method based on the "constrained half-vector symmetry" to work with non-uniform lightings. Different from these traditional methods, our method can deal with surfaces with general and unknown isotropic reflectance without the need of explicitly utilizing any additional clues or reference objects,

solving a complex optimization problem at test time, or making assumptions on the light source distribution. The work most related to ours is the UPS-FCN introduced in [5]. UPS-FCN is a single-stage model that directly regresses surface normals from images that are normalized by the known light intensities. Its performance lags far behind the calibrated methods. In contrast, our method solves the problem in two stages. We first tackles an easier problem of estimating the light directions and intensities, and then estimates the surface normals using the estimated lightings and the input images.

**Learning based lighting estimation** Recently, learning based single-image lighting estimation methods have attracted considerable attention. Gardner *et al.* [10] introduced a CNN for estimating HDR environment lighting from an indoor scene image. Hold-Geoffroy *et al.* [14] learned outdoor lighting using a physically-based sky model. Weber *et al.* [32] estimated indoor environment lighting from an image of an object with known shape. Zhou *et al.* [37] estimated lighting, in the form of Spherical Harmonics, from a human face image by assuming a Lambertian reflectance model. Different from the above methods, our method can estimate accurate directional lightings from multiple images of a static object with general shape and non-Lambertian surface.

## 3. Image Formation Model

Following the conventional practice, we assume an orthographic camera with linear radiometric response, white directional lightings coming from the upper-hemisphere, and the viewing direction pointing towards the viewer. In the rest of this paper, we refer to light direction and intensity as "lighting". Consider a non-Lambertian surface whose appearance is described by a general isotropic BRDF $\rho$. Given a surface point with normal $\boldsymbol{n} \in \mathbb{R}^3$ being illuminated by the $j$-th incoming lighting with direction $\boldsymbol{l}_j \in \mathbb{R}^3$ and intensity $e_j \in \mathbb{R}$, the image formation model can be expressed as

$$m_j = e_j \rho(\boldsymbol{n}, \boldsymbol{l}_j) \max(\boldsymbol{n}^\top \boldsymbol{l}_j, 0) + \epsilon_j, \qquad (1)$$

where $m$ represents the measured intensity, $\max(:, 0)$ accounts for attached shadows, and $\epsilon$ accounts for the global illumination effects (cast shadows and inter-reflections) and noise.

Based on this model, given the observations of $p$ surface points under $q$ different incoming lightings, the goal of uncalibrated photometric stereo is to estimate the surface normals for these $p$ surface points given only the measured intensities. In this work, we tackle this problem using a two-stage approach. In particular, we first estimate lightings from the measured intensities, and then solve for the surface normals using the estimated lightings and measured intensities.
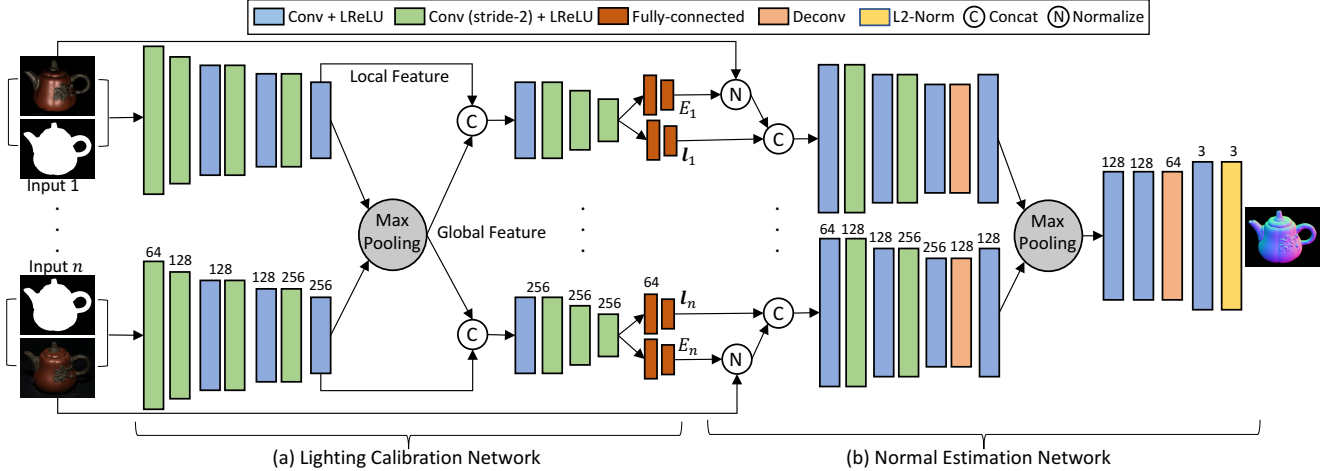
Figure 1. The network architecture of SDPS-Net is composed of (a) Lighting Calibration Network and (b) Normal Estimation Network. Kernel sizes for all convolutional layers are $3 \times 3$, and values above the layers indicate the number of feature channels.

# 4. Learning Uncalibrated Photometric Stereo

In this section, we introduce our two-stage framework, called SDPS-Net, for uncalibrated photometric stereo (see Fig. 1). The first stage of SDPS-Net, denoted as *Lighting Calibration Network* (LCNet, Fig. 1 (a)), takes an arbitrary number of images as input and estimates their corresponding light directions and intensities. The second stage of SDPS-Net, denoted as *Normal Estimation Network* (NENet, Fig. 1 (b)), estimates an accurate normal map of the object based on the lightings estimated by LCNet and the input images.

## 4.1. Lighting Calibration Network

To estimate lightings from the images, an intuitive approach would be directly regressing the light direction vectors and intensity values. However, we propose that formulating the lighting estimation as a classification problem is a superior choice, as will be verified by our experiments. Our arguments are as follows. Fist, classifying a light direction into a certain range is easier than regressing the exact value(s), and this will reduce the learning difficulty. Second, taking discretized light directions as input may allow NENet to better tolerate small errors in the estimated light directions.

**Discretization of lighting space** Since we cast our lighting estimation as a classification problem, we need to discretize the continuous lighting space. Note that a light direction in the upper-hemisphere can be described by its azimuth $\phi \in [0°, 180°]$ and elevation $\theta \in [-90°, 90°]$ (see Fig. 2 (a)). We can discretize the light directoin space by evenly dividing both the azimuth and elevation into $K_d$ bins, resulting in $K_d^2$ classes (see Fig. 2 (b)). Solving a $K_d^2$-class classification problem is not computationally efficient,
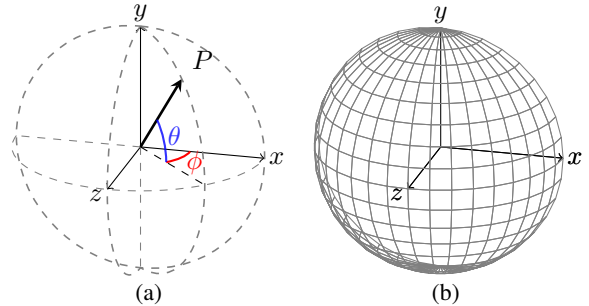


Figure 2. (a) Illustration of the coordinate system ($z$ axis is the viewing direction). $\phi \in [0°, 180°]$ and $\theta \in [-90°, 90°]$ are the azimuth and elevation of the light direction, respectively. (b) Example discretization of the light direction space when $K_d = 18$.

as the softmax probability vector will have a very high dimension even when $K_d$ is not large (*e.g.*, $K_d^2 = 1,296$ when $K_d = 36$). Instead, we estimate the azimuth and elevation of a light direction separately, leading to two $K_d$-class classification problems. Similarly, we evenly divide the range of possible light intensities into $K_e$ classes (*e.g.*, $K_e = 20$ for a possible light intensity range of $[0.2, 2.0]$).

**Local-global feature fusion** A straightforward approach to estimate the lighting for each image is simply taking a single image as input, encoding it into a feature map using a CNN, and feeding the feature map to a lighting prediction layer. It is not surprising that the result of such a simple solution is far from satisfactory. Note that the appearance of an object is determined by its surface geometry, reflectance model and the lighting. The feature map extracted from a single observation obviously does not provide sufficient information for resolving the shape-light ambiguity. Thanks to the nature of photometric stereo where multiple observations of an object are considered, we propose a local-global

feature fusion strategy to extract more comprehensive information from multiple observations.

Specifically, we separately feed each image into a shared-weight feature extractor to extract a feature map, which we call *local feature* as it only provides information from a single observation. All local features of the input images are then aggregated into a *global feature* through a max-pooling operation, which has been proven to be efficient and robust on aggregating salient features from a varying number of unordered inputs [33, 5]. Such a global feature is expected to convey implicit surface geometry and reflectance information of the object which help resolve the ambiguity in lighting estimation. Each local feature is concatenated with the global feature, and fed to a shared-weight lighting estimation sub-network to predict the lighting for each individual image. By taking both local and global features into account, our model can produce much more reliable results than using the local features alone. We empirically found that additionally including the object mask as input can effectively improve the performance of lighting estimation, as will be seen in the experiment section.

**Network architecture** LCNet is a multi-input-multi-output (MIMO) network that consists of a shared-weight *feature extractor*, an *aggregation layer* (*i.e.*, max-pooling layer), and a shared-weight *lighting estimation sub-network* (see Fig. 1 (a)). It takes the observations of the object together with the object mask as input, and outputs the light directions and intensities in the form of softmax probability vectors of dimension $K_d$ (azimuth), $K_d$ (elevation) and $K_e$ (intensity), respectively. We convert the output of LCNet to 3-vector light directions and scalar intensity values by simply taking the middle value of the range with the highest probability[1].

**Loss function** Multi-class cross entropy loss is adopted for both light direction and intensity estimation, and the overall loss function is

$$\mathcal{L}_{\text{Light}} = \lambda_{l_a}\mathcal{L}_{l_a} + \lambda_{l_e}\mathcal{L}_{l_e} + \lambda_e\mathcal{L}_e, \tag{2}$$

where $\mathcal{L}_{l_a}$ and $\mathcal{L}_{l_e}$ are the loss terms for azimuth and elevation of the light direction, and $\mathcal{L}_e$ is the loss term for light intensity. During training, weights $\lambda_{l_a}$, $\lambda_{l_e}$ and $\lambda_e$ for the loss terms are set to 1.

### 4.2. Normal Estimation Network

NENet is a multi-input-single-output (MISO) network. The network architecture of NENet is similar to PS-FCN [5], consisting of a shared-weight *feature extractor*, an *aggregation layer*, and a *normal regression sub-network* (see Fig. 1 (b)). The key difference between NENet and

PS-FCN is that PS-FCN requires accurate lightings as input, whereas NENet is trained with discretized lightings estimated by the LCNet and shows a more robust behavior over noise in the lightings.

NENet first normalizes the input images using the light intensities predicted by LCNet, and then concatenates the light directions predicted by LCNet with the images to form the input of the shared-weight feature extractor. Given an image of size $h \times w$, the loss function for NENet is

$$\mathcal{L}_{\text{Normal}} = \frac{1}{hw}\sum_i^{hw}\left(1 - \boldsymbol{n}_i^\top\tilde{\boldsymbol{n}}_i\right), \tag{3}$$

where $\boldsymbol{n}_i$ and $\tilde{\boldsymbol{n}}_i$ denote the predicted normal and the ground-truth normal, respectively, at pixel $i$.

### 4.3. Training Data

We adopted the publicly available synthetic Blobby and Sculpture datasets [5] for training. Blobby and Sculpture datasets provide surfaces with complex normal distributions and diverse materials from MERL dataset [21]. Effects of cast shadow and inter-reflection were considered during rendering using the physically based raytracer Mitsuba [16]. There are $85,212$ samples in total. Each sample was rendered under $64$ distinct light directions sampled from the upper-hemisphere with uniform light intensity, resulting in $5,453,568$ images ($85,212 \times 64$). The rendered images have a dimension of $128 \times 128$.

To simulate images under different light intensities, we randomly generated light intensities in the range of $[0.2, 2.0]$ to scale the magnitude of the images (*i.e.*, the ratio of the highest light intensity to the lowest one is $10$)[2]. Note that this selected range contains a wider range of intensity value than the public photometric stereo datasets like DiLiGenT benchmark [28] and Gourd&Apple dataset [1]. The color intensities of the input images were normalized to the range of $[0, 1]$. During training, we applied noise perturbation in the range of $[-0.025, 0.025]$ for data augmentation, and the input image size for LCNet and NENet was $128 \times 128$. At test time, NENet can take images of different dimensions, while the input for LCNet is rescaled to $128 \times 128$ as it contains fully-connected layers and requires the input to have a fixed spatial dimension. Trained only on the synthetic dataset, we will show that our model can generalize well on real datasets.

## 5. Experimental Results

We performed network analysis for our method, and compared our method with the previous state-of-the-art methods on both synthetic and real datasets.

---

[1]We have experimentally verified that alternative ways like taking the expectation of the probability vector or performing quadratic interpolation in the neighborhood of the peak value do not improve the result.

[2]Note that the ratio (other than the exact value) matters, since light intensity can only be estimated up to a scale factor.

**Implementation details** Our framework was implemented in PyTorch [24] and Adam optimizer [17] was used with default parameters. LCNet and NENet contain 4.4 million and 2.2 million parameters, respectively. We first trained LCNet using a batch size of 32 for 20 epochs until convergence, and then trained NENet from scratch given the lightings estimated by LCNet with a batch size of 16 for 10 epochs. We found that end-to-end fine-tuning did not improve the performance. The learning rate was initially set to 0.0005 and halved every 5 and 2 epochs for LCNet and NENet, respectively. It took about 22 hours to train LCNet and 26 hours to train NENet on a single Titan X Pascal GPU with a fixed input image number of 32.

**Evaluation metrics** To measure the accuracy of the predicted light directions and surface normals, the widely used mean angular error (MAE) in degree is adopted. Since the light intensities among the testing images can only be estimated up to a scale factor $s$, we introduce the scale-invariant relative error

$$E_{err} = \frac{1}{q} \sum_i^q \left( \frac{|se_i - \tilde{e}_i|}{\tilde{e}_i} \right), \qquad (4)$$

where $q$ is the number of images, $e_i$ and $\tilde{e}_i$ are the estimated and ground-truth light intensities, respectively, for image $i$. The scale factor $s$ is computed by solving $\underset{s}{\mathrm{argmin}} \sum_i^n (se_i - \tilde{e}_i)^2$ with least squares[3].

## 5.1. Network Analysis with Synthetic Data

**MERL$^{Test}$ dataset** To quantitatively perform network analysis for our method, we rendered a synthetic dataset, denoted as MERL$^{Test}$, of sphere and bunny shapes, denoted as SPHERE and BUNNY hereafter, respectively, using the physically based raytracer Mitsuba [16]. Each shape was rendered with 100 isotropic BRDFs from MERL dataset [21] under 100 light directions sampled from the upper-hemisphere, leading to 200 test objects (see Fig. 3). Cast shadows and inter-reflections are considered for BUNNY. For all experiments on synthetic dataset involving input with unknown light intensities, we randomly generated light intensities in the range of $[0.2, 2.0]$. Each experiment was repeated five times and the average results were reported.

**Discretization of lighting space** For a given number of bins $K_d$, the maximum deviation angle for azimuth and elevation of a light direction is $\delta = 180°/(K_d \times 2)$ after discretization (*e.g.*, $\delta = 2.5°$ when $K_d = 36$). To investigate how light direction discretization affects the normal estimation accuracy, we adopted the state-of-the-art calibrated method PS-FCN [5] and MERL$^{Test}$ dataset as the testbed.

---

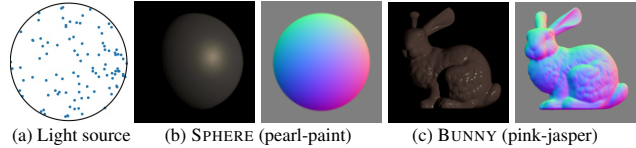(a) Light source    (b) SPHERE (pearl-paint)    (c) BUNNY (pink-jasper)

Figure 3. (a) Lighting distribution of MERL$^{Test}$ dataset. The light direction is visualized by mapping a 3-d vector $[x, y, z]$ to a point $[x, y]$. (b) and (c) show a sample image and ground-truth normal for SPHERE and BUNNY, respectively.
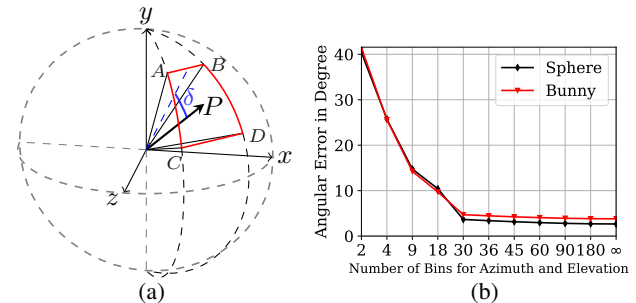


(a)          (b)

Figure 4. (a) Light directions $A, B, C$, and $D$ have the maximum deviation angles with the light direction $P$ after discretization. (b) Upper-bound of normal estimation error for PS-FCN [5] under different light direction space discretization levels ($\infty$ indicates no discretization).

We divided the azimuth and elevation of light directions into different number of bins ranging from 2 to 180. For a specific bin number, we replaced each ground-truth light direction by each of the four light directions having the maximum possible angular deviations after discretization (see Fig. 4 (a)), respectively. We then used those light directions as input for PS-FCN to infer surface normals. The normal estimation error reported in Fig. 4 (b) is the upper-bound error for PS-FCN caused by the discretization. We can see that the increase in error caused by discretization is marginal when $K_d \geq 30$. In our implementation, we empirically set $K_d$ and $K_e$ to 36 and 20, respectively. We experimentally found that the performance of LCNet is robust to different discretization levels. We chose a relatively sparse discretization of lighting space in this paper as it may allow NENet to learn to better tolerate small errors in the estimated lighting at test time.

**Effectiveness of LCNet** To validate the design of LCNet, we compared LCNet with three baseline models for lighting estimation. The first baseline model, denoted as LCNet$_{reg}$, is a regression based model that directly regresses the light direction vectors and intensity values (please refer to the supplementary for implementation details). The second baseline model, denoted as LCNet$_{w/o\ mask}$, is a classification based model that only takes the images as input without the object mask input. The last baseline model, denoted as LCNet$_{local}$, is a classification based model that independently estimates lighting for each observation (*i.e.*, without

Table 1. Lighting estimation results on the MERL$^{\text{Test}}$ dataset. The results are averaged over samples rendered with 100 BRDFs.

| ID | Model | SPHERE | | BUNNY | |
|---|---|---|---|---|---|
| | | Direction | Intensity | Direction | Intensity |
| A0 | LCNet | **3.47** | **0.082** | **5.38** | **0.089** |
| A1 | LCNet$_{\text{reg}}$ | 4.10 | 0.104 | 5.46 | 0.094 |
| A2 | LCNet$_{\text{w/o mask}}$ | 5.46 | 0.104 | 8.85 | 0.144 |
| A3 | LCNet$_{\text{local}}$ | 6.87 | 0.198 | 9.98 | 0.255 |

local-global feature fusion). All models were trained under the same setting, and the results are summarized in Table 1.

Experiments with IDs A0 & A1 in Table 1 show that the proposed classification based LCNet consistently outperformed the regression based baseline on both light direction and intensity estimation. This echoes our hypothesis that classifying a light direction to a certain range is easier than regressing an exact value. Thus, solving the classification problem reduces the learning difficulty and improves the performance. Experiments with IDs A0 & A2 show that taking the object mask as input can effectively improve the lighting estimation results. This might be explained by the fact that object mask provides strong information for occluding contours of the object, and helps the network distinguish the shadow region from the non-object region. Experiments with IDs A0 & A3 show that the proposed local-global feature fusion strategy can effectively make use of information from multiple observations, and significantly improve the lighting estimation accuracy. Please refer to our supplementary for detailed lighting estimation results of LCNet on BUNNY from MERL$^{\text{Test}}$ dataset.

**Effectiveness of NENet**  Experiments with IDs B1 & B2 in Table 2 show that after training with the discretized lightings estimated by LCNet, NENet performs better than PS-FCN given possibly noisy lightings at test time, while experiments with IDs B3 & B4 show that training NENet with the light directions estimated by the regression based baseline is not always helpful. This result further demonstrates that the proposed framework is robust to noisy lightings. Experiments with IDs B0 & B1 show the proposed method achieved results comparable to the fully calibrated method PS-FCN [5], with average MAEs of 2.71 and 4.09 on SPHERE and BUNNY, respectively.

Figure 5 shows that the performances of LCNet and NENet increased with the number of input images. This is expected, since more useful information can be used to infer the lightings and normals with more input images.

**Comparison with single-stage models**  To validate the effectiveness of the proposed two-stage framework, we compared our method with five different single-stage baseline models. We first retrained UPS-FCN [5], denoted as UPS-FCN$_{\text{retrain}}$, with images scaled by randomly generated light intensities to allow it adapt to unknown intensities at test time. We then increased the model capacity of UPS-

Table 2. Normal estimation results on the MERL$^{\text{Test}}$ dataset. The numbers are the average MAE over samples rendered with 100 BRDFs (value the lower the better). NENet$^{\dagger}$ was trained given the lightings estimated by LCNet$_{\text{reg}}$.

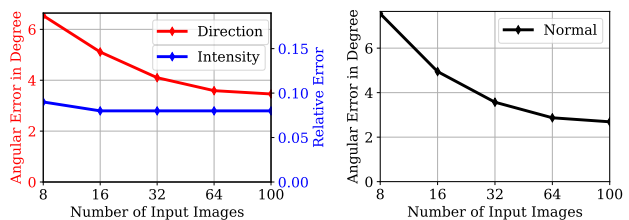| ID | Model | # Param | SPHERE | BUNNY |
|---|---|---|---|---|
| B0 | PS-FCN [5] | 2.2 M | 2.66 | 3.80 |
| B1 | LCNet + NENet | 6.6 M | **2.71** | **4.09** |
| B2 | LCNet + PS-FCN | 6.6 M | 3.19 | 4.67 |
| B3 | LCNet$_{\text{reg}}$ + NENet$^{\dagger}$ | 6.6 M | 3.22 | 4.99 |
| B4 | LCNet$_{\text{reg}}$ + PS-FCN | 6.6 M | 3.73 | 4.96 |
| B5 | UPS-FCN$_{\text{deep+mask}}$ | 6.1 M | 3.65 | 6.41 |
| B6 | UPS-FCN$_{\text{deep}}$ | 6.1 M | 4.30 | 7.29 |
| B7 | UPS-FCN$_{\text{wide}}$ | 6.4 M | 5.61 | 8.85 |
| B8 | UPS-FCN$_{\text{est\_light}}$ | 5.7 M | 6.80 | 10.62 |
| B9 | UPS-FCN$_{\text{retrain}}$ | 2.2 M | 7.44 | 12.34 |



Figure 5. Results of SDPS-Net on SPHERE from MERL$^{\text{Test}}$ dataset with varying input image numbers.

FCN by introducing a wider network (*i.e.*, more channels in the convolutional layers) and a deeper network (*i.e.*, more convolutional layers), denoted as UPS-FCN$_{\text{wide}}$ and UPS-FCN$_{\text{deep}}$, respectively. We also trained a deeper network, denoted as UPS-FCN$_{\text{deep+mask}}$, that takes both the images and object mask as input. We last investigated the effect of having additional lighting supervision by training a variant model, denoted as UPS-FCN$_{\text{est\_light}}$, to simultaneously estimate lighting and surface normal. Please refer to our supplementary for detailed network architectures.

Experiments with IDs B5-B9 in Table 2 show that utilizing a wider or deeper network, taking the object mask as input, or incorporating additional lighting supervision can improve the performance of single-stage model in some extent. However, experiments with IDs B1 & B5 show that the proposed method significantly outperformed the best-performing single-stage model, especially on surfaces with complex geometry such as BUNNY, when the input as well as the number of parameters are comparable. This result indicates that simply increasing the layer numbers or channel numbers of the network, or incorporating additional lighting supervision cannot produce optimal results.

**Comparison with the non-learning method [23]**  To further verify the effectiveness of our method over non-learning method, we compared SDPS-Net with the existing uncalibrated method PF14 [23], which achieved state-of-the-art results on the DiLiGenT benchmark [28], on differ-

Near uniform  Biased  Lambertian  Fabric  Plastic  Phenolic

(a) Light sources        (b) Examples for the four typical types of BRDFs

(c) Light direction estimation results

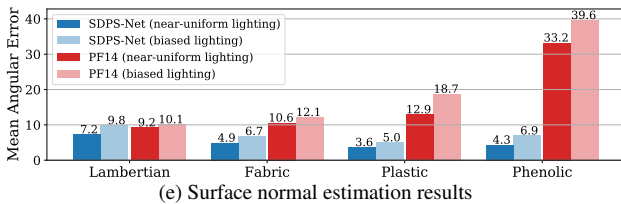(d) Light intensity estimation results
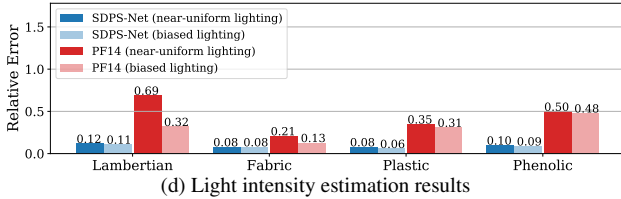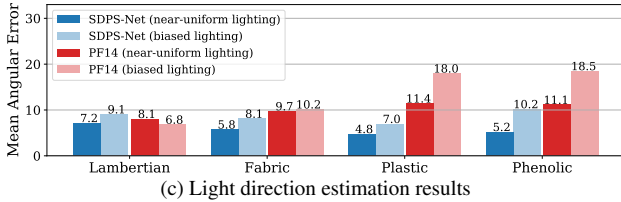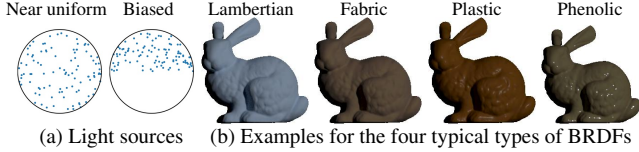
(e) Surface normal estimation results

Figure 6. Comparison between SDPS-Net and PF14 [23] on BUNNY rendered with four different types of BRDFs under a near uniform lighting distribution and a biased lighting distribution.

ent lighting distributions and types of BRDFs. Specifically, we considered one near uniform and one biased lighting distribution (see Fig. 6 (a)). We rendered BUNNY using four typical types of BRDFs, including the Lambertian model and three other types from MERL dataset [21], namely, Fabric, Plastic, and Phenolic. They contained 15, 12, 9, and 12 different BRDFs, respectively. We reported the average results for each type (see Fig. 6 (b) for an example of each type.).

Figures 6 (c)-(e) compare SDPS-Net and PF14 on lighting estimation and normal estimation. The following observations are made: 1) PF14 performed well on light direction and normal estimation for diffuse or near diffuse surfaces (*i.e.*, Lambertian and Fabric), but will quickly degenerate when dealing with non-Lambertian surfaces. Besides, it cannot reliably estimate light intensities for all the BRDFs. 2) SDPS-Net performed well on different types of BRDFs, especially on surfaces exhibit specular highlights. This result suggests that specular highlight is an important clue for uncalibrated photometric stereo [7]. 3) The performance of light direction and normal estimation of both methods will have a trend of decreasing when dealing with biased lighting distribution, while the performance of intensity estimation will slightly improve.



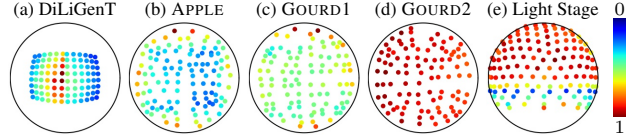(a) DiLiGenT  (b) APPLE  (c) GOURD1  (d) GOURD2  (e) Light Stage

Figure 7. Lighting distributions of real testing datasets. The light direction is visualized by mapping a 3-d vector $[x, y, z]$ to a point $[x, y]$. The color of the point indicates the light intensity (value is divided by the highest intensity to normalize to $[0, 1]$).

## 5.2. Evaluation on Real Datasets

**Real testing datasets** We evaluated our method on three publicly available non-Lambertian photometric stereo datasets, namely the *DiLiGenT benchmark* [28], *Gourd&Apple dataset* [1] and *Light Stage Data Gallery* [8]. Figure 7 visualizes the lighting distribution of these datasets (note that for Light Stage Data Gallery, we only used 133 images with the front side of the object under illumination). Since Gourd&Apple dataset and Light Stage Data Gallery only provide calibrated lightings (without ground-truth normal maps), we quantitatively evaluated our method on lighting estimation while qualitatively evaluated it on normal estimation.

**Evaluation on DiLiGenT benchmark** Table 3 (a)-(b) show that LCNet outperformed the regression based baseline LCNet$_{reg}$ and achieved highly accurate results on both light direction and intensity estimation on DiLiGenT benchmark, with an average MAE of 4.92 and an average relative error of 0.068, respectively. Table 3 (c) compares the normal estimation results of SDPS-Net with previous state-of-the-art methods on DiLiGenT benchmark. SDPS-Net achieved state-of-the-art results on almost all objects with an average MAE of 9.51, except for the BEAR object. Although UPS-FCN$_{deep+mask}$ achieved reasonably good results on objects with smooth surface and uniform material (*e.g.*, BALL), it had difficulties in handling surfaces with complex geometry and spatially-varying BRDFs (*e.g.*, READING and HARVEST). The normal estimation network coupled with LCNet (i.e., SDPS-Net) outperforms that with LCNet$_{reg}$ (i.e., LCNet$_{reg}$+NENet$^\dagger$) with a clear improvement of 1.52 in average MAE, demonstrating the effectiveness of the proposed classification based LCNet. It is interesting to see that, coupled with our LCNet, the calibrated methods L2 baseline [34] and IS18 [15] can already achieve results comparable to the previous state-of-the-art methods. This result indicates that our proposed LCNet can be integrated with existing calibrated methods to help handle cases where lighting conditions are unknown. Figures 8 (a)-(b) show the qualitative results of SDPS-Net on DiLiGenT benchmark.

**Evaluation on other real datasets** Table 4 shows that SDPS-Net can estimate accurate light directions and intensities for the challenging Gourd&Apple dataset and Light
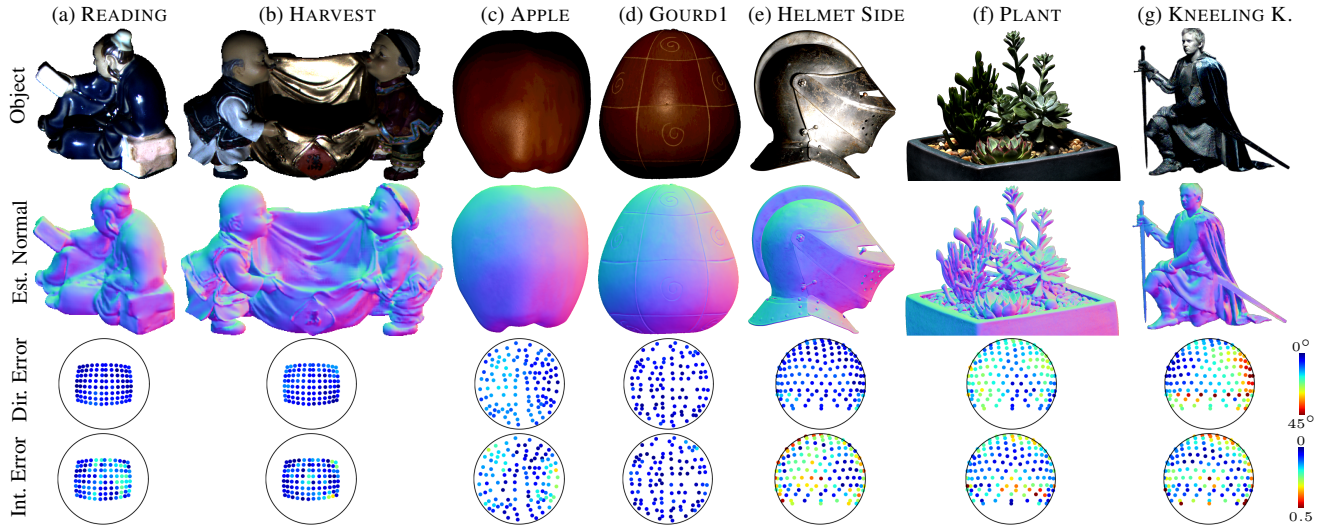
Figure 8. Qualitative results of SDPS-Net on the real testing datasets. The first to the fourth rows show the object, estimated normal map, error distribution of light direction and light intensity estimation, respectively.

Table 3. Results of SDPS-Net on the DiLiGenT benchmark.

(a) Results on light direction estimation.

| Method | BALL | CAT | POT1 | BEAR | POT2 | BUDDHA | GOBLET | READING | COW | HARVEST | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|
| LCNet$_{reg}$ | 4.94 | 5.82 | 5.62 | 7.19 | 4.82 | **3.90** | 12.89 | 7.90 | **4.19** | 9.50 | 6.68 |
| LCNet | **3.27** | **4.08** | **5.44** | **3.47** | **2.87** | 4.34 | **10.36** | **4.50** | 4.52 | **6.32** | **4.92** |

(b) Results on light intensity estimation.

| Method | BALL | CAT | POT1 | BEAR | POT2 | BUDDHA | GOBLET | READING | COW | HARVEST | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|
| LCNet$_{reg}$ | **0.032** | **0.051** | **0.048** | 0.167 | 0.074 | 0.080 | 0.075 | 0.141 | **0.044** | 0.085 | 0.080 |
| LCNet | 0.039 | 0.095 | 0.058 | **0.061** | **0.048** | **0.048** | **0.067** | **0.105** | 0.073 | **0.082** | **0.068** |

(c) Results on normal estimation. (Best viewed in PDF with zoom.)

| Method | BALL | CAT | POT1 | BEAR | POT2 | BUDDHA | GOBLET | READING | COW | HARVEST | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|
| AM07 [2] | 7.27 | 31.45 | 18.37 | 16.81 | 49.16 | 32.81 | 46.54 | 53.65 | 54.72 | 61.70 | 37.25 |
| SM10 [27] | 8.90 | 19.84 | 16.68 | 11.98 | 50.68 | 15.54 | 48.79 | 26.93 | 22.73 | 73.86 | 29.59 |
| WT13 [35] | 4.39 | 36.55 | 9.39 | **6.42** | 14.52 | 13.19 | 20.57 | 58.96 | 19.75 | 55.51 | 23.93 |
| LM13 [19] | 22.43 | 25.01 | 32.82 | 15.44 | 20.57 | 25.76 | 29.16 | 48.16 | 22.53 | 34.45 | 27.63 |
| PF14 [23] | 4.77 | 9.54 | 9.51 | 9.07 | 15.90 | 14.92 | 29.93 | 24.18 | 19.53 | 29.21 | 16.66 |
| LC18 [18] | 9.30 | 12.60 | 12.40 | 10.90 | 15.70 | 19.00 | 18.30 | 22.30 | 15.00 | 28.00 | 16.30 |
| UPS-FCN [5] | 6.62 | 14.68 | 13.98 | 11.23 | 14.19 | 15.87 | 20.72 | 23.26 | 11.91 | 27.79 | 16.02 |
| LCNet + L2 [34] | 4.90 | 11.12 | 9.72 | 9.35 | 14.70 | 14.86 | 18.29 | 20.11 | 25.08 | 29.17 | 15.73 |
| LCNet + IS18 [15] | 6.37 | 15.64 | 10.58 | 8.48 | 12.24 | 13.94 | 18.54 | 23.78 | 25.69 | 16.46 |
| UPS-FCN$_{deep+mask}$ | 3.96 | 12.16 | 11.13 | 7.19 | 11.11 | 13.06 | 18.07 | 20.46 | 11.84 | 27.22 | 13.62 |
| LCNet$_{reg}$+NENet$^{\dagger}$ | 3.87 | 8.97 | **8.04** | 15.98 | 8.36 | 9.42 | **11.49** | 16.99 | 8.83 | 18.38 | 11.03 |
| SDPS-Net | **2.77** | **8.06** | 8.14 | 6.89 | **7.50** | **8.97** | 11.91 | **14.90** | 8.48 | 17.43 | 9.51 |

Table 4. Lighting estimation results of SDPS-Net on the Gourd&Apple dataset and the Light Stage Data Gallery.

(a) Results on the Gourd&Apple dataset.

| | APPLE | GOURD1 | GOURD2 | Avg. |
|---|---|---|---|---|
| Direction | 9.31 | 4.07 | 7.11 | 6.83 |
| Intensity | 0.106 | 0.048 | 0.186 | 0.113 |

(b) Results on the Light Stage Data Gallery.

| | HELMET SIDE | PLANT | FIGHTING KNIGHT | KNEELING KNIGHT | STANDING KNIGHT | HELMET FRONT | Avg. |
|---|---|---|---|---|---|---|---|
| Direction | 6.57 | 16.06 | 15.95 | 19.84 | 11.60 | 11.62 | 13.61 |
| Intensity | 0.212 | 0.170 | 0.214 | 0.199 | 0.286 | 0.248 | 0.221 |

Stage Data Gallery. Our method can also reliably recover visually pleasing surface normal of these two datasets (see Fig. 8 (c)-(g)), clearly demonstrating the practicality of the proposed methods in real world applications. Please refer to our supplementary for more results.

## 6. Conclusion and Discussion

In this paper, we have proposed a two-stage deep learning framework, called SDPS-Net, for uncalibrated photometric stereo. The first stage of our framework takes an arbitrary number of images as input and estimates their corresponding light directions and intensities, while the second stage predicts the normal map of the object based on the lightings estimated in the first stage and the input images. By explicitly learning to estimate lighting conditions, our two-stage framework can take advantage of the intermediate supervision to reduce the learning difficulty and improve the final normal estimation results. Besides, the first stage of our framework can be seamlessly integrated with existing calibrated methods, which enables them to handle uncalibrated photometric stereo. Experiments on both synthetic and real datasets showed that our method significantly outperformed existing state-of-the-art uncalibrated photometric stereo methods.

Since our framework is trained only on surfaces with uniform material, it may not perform well in dealing with steep color changes caused by multi-material surfaces (see Fig. 8 (b) for an example). In the future, we will investigate better training datasets and network architectures for handling surfaces with spatially-varying BRDFs.

# References

[1] Neil Alldrin, Todd Zickler, and David Kriegman. Photometric stereo with non-parametric and spatially-varying reflectance. In *CVPR*, 2008. 4, 7

[2] Neil G Alldrin, Satya P Mallick, and David J Kriegman. Resolving the generalized bas-relief ambiguity by entropy minimization. In *CVPR*, 2007. 1, 2, 8

[3] Peter N Belhumeur, David J Kriegman, and Alan L Yuille. The bas-relief ambiguity. *IJCV*, 1999. 1, 2

[4] Manmohan Krishna Chandraker, Fredrik Kahl, and David J Kriegman. Reflections on the generalized bas-relief ambiguity. In *CVPR*, 2005. 2

[5] Guanying Chen, Kai Han, and Kwan-Yee K. Wong. PS-FCN: A flexible learning framework for photometric stereo. In *ECCV*, 2018. 1, 2, 4, 5, 6, 8

[6] Donghyeon Cho, Yasuyuki Matsushita, Yu-Wing Tai, and Inso Kweon. Photometric stereo under non-uniform light intensities and exposures. In *ECCV*, 2016. 2

[7] Ondrej Drbohlav and M Chaniler. Can two specular pixels calibrate photometric stereo? In *ICCV*, 2005. 2, 7

[8] Per Einarsson, Charles-Felix Chabert, Andrew Jones, Wan-Chun Ma, Bruce Lamond, Tim Hawkins, Mark Bolas, Sebastian Sylwan, and Paul Debevec. Relighting human locomotion with flowed reflectance fields. In *EGSR*, 2006. 7

[9] Carlos Hernandez Esteban, George Vogiatzis, and Roberto Cipolla. Multiview photometric stereo. *IEEE TPAMI*, 2008. 2

[10] Marc-André Gardner, Kalyan Sunkavalli, Ersin Yumer, Xiaohui Shen, Emiliano Gambaretto, Christian Gagné, and Jean-François Lalonde. Learning to predict indoor illumination from a single image. *ACM TOG*, 2017. 2

[11] Athinodoros S Georghiades. Incorporating the torrance and sparrow model of reflectance in uncalibrated photometric stereo. In *ICCV*, 2003. 2

[12] Hideki Hayakawa. Photometric stereo under a light source with arbitrary motion. *JOSA A*, 1994. 2

[13] Aaron Hertzmann and Steven M Seitz. Example-based photometric stereo: Shape reconstruction with general, varying brdfs. *IEEE TPAMI*, 2005. 2

[14] Yannick Hold-Geoffroy, Kalyan Sunkavalli, Sunil Hadap, Emiliano Gambaretto, and Jean-François Lalonde. Deep outdoor illumination estimation. In *CVPR*, 2017. 2

[15] Satoshi Ikehata. CNN-PS: CNN-based photometric stereo for general non-convex surfaces. In *ECCV*, 2018. 1, 2, 7, 8

[16] Wenzel Jakob. Mitsuba renderer, 2010. 4, 5

[17] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 5

[18] Feng Lu, Xiaowu Chen, Imari Sato, and Yoichi Sato. Symps: BRDF symmetry guided photometric stereo for shape and light source estimation. *IEEE TPAMI*, 2018. 2, 8

[19] Feng Lu, Yasuyuki Matsushita, Imari Sato, Takahiro Okabe, and Yoichi Sato. Uncalibrated photometric stereo for unknown isotropic reflectances. In *CVPR*, 2013. 1, 2, 8

[20] Feng Lu, Imari Sato, and Yoichi Sato. Uncalibrated photometric stereo based on elevation angle recovery from BRDF symmetry of isotropic materials. In *CVPR*, 2015. 1

[21] Wojciech Matusik, Hanspeter Pfister, Matt Brand, and Leonard McMillan. A data-driven reflectance model. In *SIGGRAPH*, 2003. 4, 5, 7

[22] Takahiro Okabe, Imari Sato, and Yoichi Sato. Attached shadow coding: Estimating surface normals from shadows under unknown reflectance and lighting conditions. In *ICCV*, 2009. 2

[23] Thoma Papadhimitri and Paolo Favaro. A closed-form, consistent and robust solution to uncalibrated photometric stereo via local diffuse reflectance maxima. *IJCV*, 2014. 1, 2, 6, 7, 8

[24] Adam Paszke, Sam Gross, Soumith Chintala, and Gregory Chanan. PyTorch: Tensors and dynamic neural networks in python with strong gpu acceleration, 2017. 5

[25] Hiroaki Santo, Masaki Samejima, Yusuke Sugano, Boxin Shi, and Yasuyuki Matsushita. Deep photometric stereo network. In *ICCV Workshops*, 2017. 1, 2

[26] Imari Sato, Takahiro Okabe, Qiong Yu, and Yoichi Sato. Shape reconstruction based on similarity in radiance changes under varying illumination. In *ICCV*, 2007. 2

[27] Boxin Shi, Yasuyuki Matsushita, Yichen Wei, Chao Xu, and Ping Tan. Self-calibrating photometric stereo. In *CVPR*, 2010. 1, 2, 8

[28] Boxin Shi, Zhipeng Mo, Zhe Wu, Dinglong Duan, Sai-Kit Yeung, and Ping Tan. A benchmark dataset and evaluation for non-Lambertian and uncalibrated photometric stereo. *IEEE TPAMI*, 2018. 1, 2, 4, 7

[29] William M Silver. *Determining shape and reflectance using multiple images*. PhD thesis, Massachusetts Institute of Technology, 1980. 1

[30] Ping Tan, Satya P Mallick, Long Quan, David J Kriegman, and Todd Zickler. Isotropy, reciprocity and the generalized bas-relief ambiguity. In *CVPR*, 2007. 2

[31] Tatsunori Taniai and Takanori Maehara. Neural inverse rendering for general reflectance photometric stereo. In *ICML*, 2018. 1, 2

[32] Henrique Weber, Donald Prévost, and Jean-François Lalonde. Learning to estimate indoor lighting from 3d objects. In *3DV*, 2018. 2

[33] Olivia Wiles and Andrew Zisserman. SilNet: Single-and multi-view reconstruction by learning from silhouettes. In *BMVC*, 2017. 4

[34] Robert J Woodham. Photometric method for determining surface orientation from multiple images. *Optical engineering*, 1980. 1, 7, 8

[35] Zhe Wu and Ping Tan. Calibrating photometric stereo by holistic reflectance symmetry analysis. In *CVPR*, 2013. 2, 8

[36] Alan L Yuille, Daniel Snow, Russell Epstein, and Peter N Belhumeur. Determining generative models of objects under varying illumination: Shape and albedo from multiple images using SVD and integrability. *IJCV*, 1999. 2

[37] Hao Zhou, Jin Sun, Yaser Yacoob, and David W. Jacobs. Label denoising adversarial network (LDAN) for inverse lighting of faces. In *CVPR*, 2018. 2