

A CNN Based Approach for the Point-Light Photometric Stereo Problem

Fotios Logothetis* Roberto Mecca* Ignas Budvytis[†] Roberto Cipolla[†]

October 11, 2022

Abstract

Reconstructing the 3D shape of an object using several images under different light sources is a very challenging task, especially when realistic assumptions such as light propagation and attenuation, perspective viewing geometry and specular light reflection are considered. Many of works tackling Photometric Stereo (PS) problems often relax most of the aforementioned assumptions. Especially they ignore specular reflection and global illumination effects. In this work, we propose a CNN-based approach capable of handling these realistic assumptions by leveraging recent improvements of deep neural networks for far-field Photometric Stereo and adapt them to the point light setup. We achieve this by employing an iterative procedure of point-light PS for shape estimation which has two main steps. Firstly we train a per-pixel CNN to predict surface normals from reflectance samples. Secondly, we compute the depth by integrating the normal field in order to iteratively estimate light directions and attenuation which is used to compensate the input images to compute reflectance samples for the next iteration.

Our approach significantly outperforms the state-of-the-art on the DiLiGenT real world dataset. Furthermore, in order to measure the performance of our approach for near-field point-light source PS data, we introduce LUCES the first real-world 'dataset for near-field point light source photometric Stereo' of 14 objects of different materials where the effects of point light sources and perspective viewing are a lot more significant. Our approach also outperforms the competition on this dataset as well. Data and test code are available at the project page¹.

1 Introduction

Retrieving the 3D shape of a static object from observations under varying illumination is a very challenging problem in Computer Vision under the name of Photometric Stereo (PS). PS has been used in the past for inspection tasks such as the examination of the fracture of sandstone samples [26] or defect detection of steel components manufacturing [48].

Originally, [62] proposed a mathematical model of the PS problem relying on four main assumptions: orthographic viewing geometry, diffuse light reflection, uniform light propagation and the lack of global illumination effects (cast shadows, self reflections, ambient light). Due to restrictive assumptions, such method was limited to very narrowly specified scenarios. Since then an extensive research has been carried out to relax these assumptions.

Shape reconstruction from shading information is a difficult problem, due to the complexity of the underlying physical process describing how a light beam bounces on the surface. Thus, it becomes very important to take into account the parametrization of all elements that influence the image formation. After [62], most of the literature dealing with PS still assume diffuse reflection (i.e., uniform in all directions), reducing the mathematical model to a linear problem where the normal field can be easily computed [57] and finally integrated [16, 45]. Realistically, this approach contains too many assumptions which fail as soon as the method is used in a real-world application. There are at least two reasons why the reconstruction in particular of specular surfaces still remains a challenging task in

*Toshiba Europe Ltd flogothetis,rmecca@crl.toshiba.co.uk

[†]University of Cambridge ib255,rc10001@cam.ac.uk

¹<https://www.toshiba.eu/pages/eu/Cambridge-Research-Laboratory/luces>



Figure 1: Our proposed approach accurately reconstructs highly specular objects, in various datasets including DiLiGenT [51] and LUCES [37].

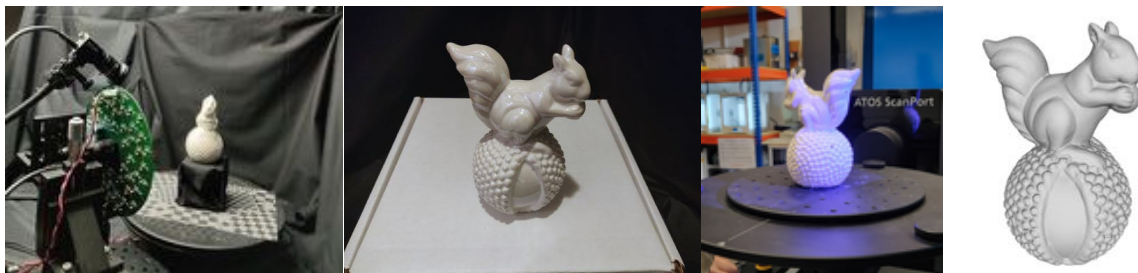


Figure 2: From left to right: (1) the stage of our Photometric Stereo setup, (2) a top view of a sample object (Squirrel), (3) acquisition with the GOM scanner, (4) the 3D scanned mesh.

the PS field. First, the Bidirectional Reflectance Distribution Function (BRDF) for specular reflections is highly non-linear, which means that analytical solutions are intractable. Second, the behaviour of light before and after it starts bouncing on the object needs to be modeled accurately. As many recent methods [21,31,50] aim at retrieving the geometry at a per-pixel level, light-object interaction requires proper modeling.

With the aim of solving the PS problem under more realistic conditions, researchers have modelled perspective viewing geometry [41,43,57], specular light reflections [38] and point light sources parameterising radial propagation of light [13,24]. These effects lead to highly non-linear models requiring sophisticated optimisation strategies [44,60]. As the complexity of the models becomes intractable, especially when dealing with physical material properties, several recent works have opted to neglect specular highlights and instead rely on robust optimisation techniques [22,23]. Furthermore, real objects experience a number of complex physical effects which make the explicit mathematical modeling very hard to invert.

In fact, these global illumination effects (cast shadows, self reflections, ambient light) are one of the most challenging aspects of PS. [34,66] tackle the case of fixed ambient light which however is too simple of a model to cover realistic inter-reflections. The global illumination issue is firstly adequately addressed in [21] by employing a Convolutional Neural Network (CNN). This method works by arranging reflectance samples into a fixed size observation map for each pixel. Observational maps are then provided as an input to the CNN which is trained to output a surface normal per pixel. [31] extends this work and shows how a training data augmentation strategy can be used to deal with

general BRDFs such as MERL dataset [35] or Bidirectional Scattering Distribution Functions (BSDF) such as Disney [8] and global illumination effects in the far-field setting. However, these approaches are only directly applicable to the far-field photometric stereo since the nonlinear light attenuation of near-field images does not allow to directly compute valid observation maps.

Usually the concept of near-field PS is relevant when the images are acquired with the camera/light setup nearby the object. Differently from the far-field case where incoming light is parametrised as a uniform 3D vector, light directions at every pixel location are dependent on the geometry of the light source. Instead of dealing with general lighting models [47], most of the approaches consider a point-like light source which actually matches the widely available LED based illumination. It is important to notice that even far-field PS datasets are acquired by using point light sources [52]. In this work we use the concept of point-light based PS and, instead of constraining it in the near-field, we provide a method which is able to improve state-of-the-art also in the far-field.

To do so, we use a three step process. Firstly, the effect of the light attenuation is compensated using an estimate of the object geometry, to produce equivalent far-field reflectance samples. Secondly, a CNN is used to regress pixel normals from these samples. Finally, a numerical integration is used to update the estimate of the object geometry for the next iteration step. We evaluate our method on both artificial and real point-light image datasets. We significantly outperform competing approaches [32, 44, 50] on both types of datasets (see Figure 1 and Section 6).

We extend our previous approach [30] by making the network able to train over a general point-light distribution. We tested over a wide variety of scenarios, taking into account sparse and dense point-light distribution as well as synthetic experiments. We finally compare our method over the real-world point-light PS datasets DiLiGenT [52] and LUCES [37] to cover both far and near field setting respectively. We also extend the preliminary version of LUCES (see Figures 2, 5 and Section 4) [37] by analysing the real to synthetic gap among a variety of competing methods. In addition, we improved ground truth meshes by employing CT scanning technologies² to retrieve 3D geometry of objects made of non-diffuse materials. Under different lighting setups between LUCES and DiLiGenT, object materials, focal length and illumination density, we discuss the variation of performances for different configuration of several approaches.

The rest of this work is divided as follows. Section 2 discusses relevant work in Photometric Stereo. Section 3 provides details of our proposed method. Section 4 outlines the LUCES dataset. Sections 5 and 6 describe the experiment setup and corresponding results.

2 Related Work

In this section we provide an overview of the relevant latest improvements in PS. For a detailed, fairly recent PS survey, refer to [3].

2.1 Point-light Source PS

Differently from the classical directional-light PS, point-light source based approaches assume that the illumination spreads non-linearly with respect to the position of the light sources, thus making analytical models more complicated and harder to solve in practice.

Most of the approaches that dealt with point-light illumination were actually trying to solve specific applications mostly related to endoscopic inspection [14, 15, 42, 63]. In particular they were always trying to tackle the problem under near-field setting, which is the most obvious scenario where non-linear light behaviour and image perspective have to be addressed in order to avoid distorted geometry.

In this particular endoscopic framework, Wu *et al.* [63] studied the multi-image endoscopic problem by considering two light sources placed off the optical center. They developed an irradiance model obtained by simultaneously illuminating an object with two different light sources. They then recovered the surface by considering a single irradiance equation for the sum of Lambertian reflectance functions of the two different light sources. The use of this reflectance function results in a loss of information. In order to avoid this problem and issues related to an unknown albedo, they used a photometric calibration. Surface recovery is performed within a variational framework that involves

²www.zeiss.com/metrology

high computational complexity compared to alternative direct methods [36]. The shape from an endoscopic perspective problem solved via a photometric stereo technique using more than 2 images was first addressed by Collins and Bartoli [14]. They solved the close-range PS with with an a-priori light calibration procedure. Furthermore, they used a prior for a reflectance model learning by adding physical markers on the inspected object even when the surface was assumed to be Lambertian. In particular, their mathematical formulation is based on the usual two step procedure where an energy functional is minimized (which allows the computation of the surface derivatives), and only later is the surface recovered [4, 16, 53]. Moreover, their energy is based on the sum of Lambertian irradiance equations rather than using photometric ratios [9, 36, 58, 61] that lead to more practical problems. For example, the most important feature of photometric ratios is to obtain independence from the albedo. Parot *et al.* [42] studied the same problem by using a straightforward heuristic approach to photometric stereo. In their work, even if camera and lights are close to the inspected object, they assumed orthographic viewing geometry, with uniform and unattenuated light directions calibrated by assuming reasonable distance between the object and the camera. The discrepancy with respect to the real physics about object proximity is faced by filtering the directional gradients depending on the frequencies. They heuristically handled this by removing the lower frequencies and the DC components. Then, the resulting depth map is computed using a multigrid Poisson solver [53]. The work describes purely qualitative results in the sense that they did not represent accurate reconstructions of the environment, instead they used their method as a qualitative tool for detecting lesions.

Some works embedded the non-linearities coming from the point-light source geometry in a Partial Differential Equation (PDE)-based formulation using image ratios [39, 40]. This way allows to calculate depth directly, without the intermediate step of approximating the normal field. Also [27, 54] took advantage of image ratios in order to eliminate the dependence on the surface albedo and thus reduce the number of unknowns. Image ratios were also used in the variational framework of [38] in order to make the approach more robust to specular highlights by unifying diffuse and Blinn-Phong specular [7] reflections into a single mathematical formulation. This general variational framework is also applicable in a weakly calibrated setting [32] or even a volumetric one [33]. Recently, a LED-based approach introduced by Quéau *et al.* [44] presented a complicated variational approach based on alternating weighted least-square scheme also capable of calibrating the light brightness of the light sources. Furthermore, [29] exploited a circular LED setup to compute the relative mean distance between the camera and the object.

2.2 Deep Learning (DL) Based Approaches for PS

Computer graphics is a well understood topic and many tools capable of rendering highly non-linear irradiance equations are publicly available³ [35]. This allowed to create reliable datasets for supervised DL approaches. The potential of DL for solving the PS problem can be divided in two main advantages. Firstly, CNNs have the capability of inverting highly non-linear reflectance models comprising of numerous physically based parameters. Secondly, CNNs can be made to deal with the complicated real world imperfection (shadows, self reflections, noise) through the use of data augmentation. So far, several DL approaches have been proposed [21, 31, 49, 50]. A preliminary work by [18, 55] considered diffuse reflection only. [65] proposed a library where set of novel layers can be incorporated into a generic neural network to embed explicit models of photometric image formation. More recently, several approaches have tackled the problem of reconstructing complex objects. [49] proposed a method to find correspondences between simulated observation rendered by the MERL BRDF dataset [35] and the normal map of the target object, handling non-local effects using a dropout strategy. [25] leveraged DL to learn the information from multispectral images to get RGB based PS reconstructions. [56] proposed generating training data on the go to minimise the image re-projection error. Although this method is a training data free approach, the whole procedure is relatively slow. [28] proposed a dedicated network to account for global illumination effects for the case of single mobile image reconstruction. Recently, [11] proposed rendering patches of different surface materials in order to get training data. This method is also extended in [10] for solving the uncalibrated PS. Ikehata [21] proposed arranging all the reflectance samples of a pixel (i.e. different illumination images in the far-field setting) into a fixed size *observation* map which is used by a CNN to regress pixel normals. The CNN is essentially learning to invert the BRDF with added robustness

³www.blender.org and www.disneyanimation.com/technology/brdf.html

to global illumination effects, as training data are made with physics based rendering. In [31], this method was extended by simplifying the training procedure providing an inline per-pixel training data generation.

However, none of these DL approaches directly tackle the point-light PS problem and non-linear attenuation from point light sources together with the viewing direction dependency drastically increase the problem space exploding the training data requirements. Santo *et al.* [50] addressed this problem with a hybrid approach where the light reflected is firstly interpreted as coming from a directional light source, and then refined with a point-light model based on a near-light image formation.

In this work, we expand our method [30] which was limited to provide depth prediction for the point-light PS problem for a specific light configuration. The proposed per-pixel training procedure has been improved in order to include a much wider variety of lighting scenarios. This allows the proposed network to provide state-of-the-art predictions in a general point-light setup.

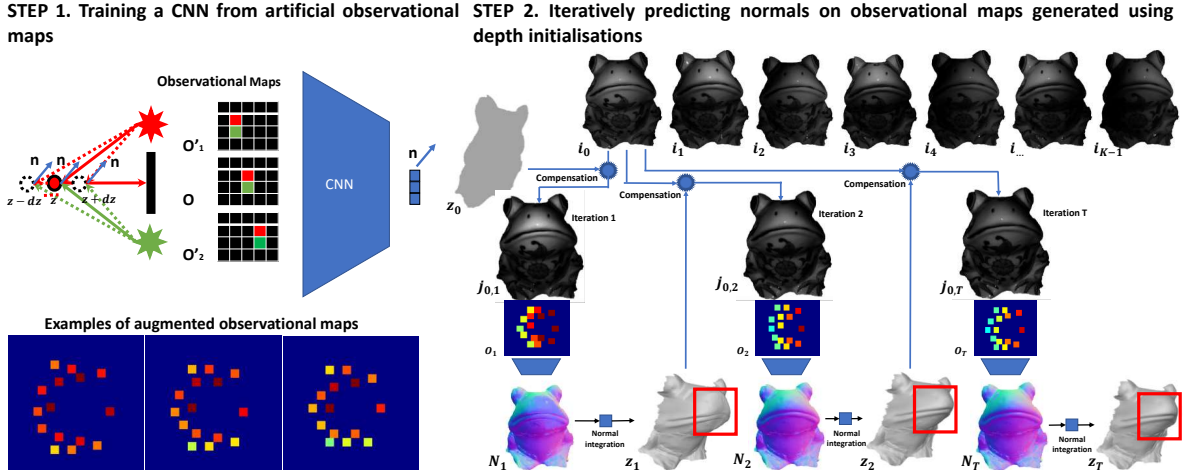


Figure 3: This figure illustrates two key steps of our proposed approach. On the left, the network training is illustrated consisting of sampling points inside the camera’s frustrum and rendering the respective observation maps. As the depth will only be approximately known at test time, this is slightly perturbed before mapping resulting to a structured change of the map (this structured change is shown at bottom left: the middle image is the map computed with the actual depth (10 cm), the left and right maps are computed with 9 and 11 cm respectively). On the right, the reconstruction process is shown. Images $i_0 \dots i_{K-1}$ are used with conjunction with previous depth estimate to compensate for light attenuation ($j_0 \dots j_{K-1}$), compute observation maps (shown for pixel at image center here), regress normals and finally update the shape. Note the improvement of the shape of frogs beak (red square) from iteration 1 to iteration 2. Also see Fig.4.

2.3 Photometric Stereo Datasets

Across the years, a number of custom real-world PS datasets have been created to suit the purposes of the proposed approaches. Alldrin *et al.* [5] proposed a dataset consisting of 3 objects lit by roughly a hundred distant directional lights. The light calibration in terms of positioning and intensity has been performed by using respectively a mirror sphere and a diffuse sphere. Xiong *et al.* [64] have proposed a dataset of 7 objects using 20 directional lights calibrated with two chrome spheres. As the approach was mostly modeling PS images with Lambertian irradiance equations, the material of the objects was quite diffuse. A limited number of PS data has been released by Quéau *et al.* to prove the working principle of an edge preserving method [45] and a multi-spectral PS approach [46].

Although initially designed for evaluating multi-view approaches, the datasets released by Aanæs *et al.* [1,2] are useful for evaluating PS approaches as they also contain images under varying illumination. As most of the methods aimed at tackling the PS problem deal with the far-field setting, recently Shi *et al.* [51] introduced the first dataset in this category, namely DiLiGenT aimed at evaluating reconstruction methods over a wide variety of materials for 10 different objects. This work also contains

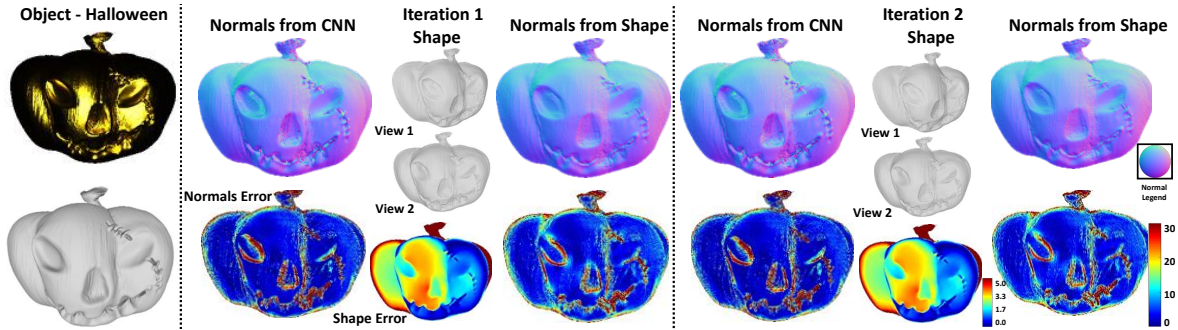


Figure 4: Iterative refinement of the geometry for the Halloween synthetic object. On the left, a sample image and GT shape are shown. The other 2 sections show 2 steps of the iterative refinement with the respective normals (both raw network predictions and differentiated ones), normal error maps and depth error maps. As the difference is minimal between steps 1 and 2, the process is converged.

a well discussed taxonomy for non-Lambertian and uncalibrated PS approaches. Their setup consists of 96 LEDs placed several meters away from the objects to approximate directional illumination and the camera (with a 50mm lens) was placed at 1.5m from the object. Such distance between the object and the camera/lights system does not provide to this dataset the near-field light variation studied in many recent approaches.

3 Method

In this section we describe our method for tackling the point-light Photometric Stereo problem. In particular, we provide both the details of the assumed image formation model and how normals can be predicted for Photometric Stereo images by using CNN’s trained on reflectance samples (also see Figure 3).

3.1 Point-light Modeling

Similar to [40], we assume calibrated point light sources at positions \mathbf{P}_m (w.r.t the camera center at $\mathbf{0}$) resulting in variable lighting vectors $\mathbf{L}_m = \mathbf{P}_m - \mathbf{X}$. Here $\mathbf{X} = [x, y, z]^T$ is the 3D surface point coordinates. We also model the light attenuation considering the following non-linear radial model of dissipation:

$$a_m(\mathbf{X}) = \phi_m \frac{(\hat{\mathbf{L}}_m(\mathbf{X}) \cdot \hat{\mathbf{D}}_m)^{\mu_m}}{\|\mathbf{L}_m(\mathbf{X})\|^2}, \quad (1)$$

where $\hat{\mathbf{L}}_m = \frac{\mathbf{L}_m}{\|\mathbf{L}_m\|}$ is the lighting direction, ϕ_m is the intrinsic brightness of the light source, $\hat{\mathbf{D}}_m$ is the principal direction (i.e. the orientation of the LED point light source) and μ_m is an angular dissipation factor. Defining $\hat{\mathbf{V}} = -\frac{\mathbf{X}}{\|\mathbf{X}\|}$ as the viewing vector, the general image irradiance equation becomes:

$$i_m = a_m \mathbf{B}(\mathbf{N}, \hat{\mathbf{L}}_m, \hat{\mathbf{V}}, \rho). \quad (2)$$

Here \mathbf{N} is the surface normal. \mathbf{B} is assumed to be a general BRDF and ρ is the surface albedo (allowing for the most general case, images and ρ are RGB and the reflectance is different per channel). In addition, we allow for the possibility of global illumination effects (shadows, self reflections) which are incorporated into \mathbf{B} . This can be re-arranged into a BRDF inversion problem as (for BRDF samples j_m):

$$j_m = \frac{i_m}{a_m} = \mathbf{B}(\mathbf{N}, \hat{\mathbf{L}}_m, \hat{\mathbf{V}}, \rho). \quad (3)$$

We note that $\hat{\mathbf{V}}$ is known but \mathbf{L}_m and a_m are unknowns due to the nonlinear dependence on z . Our objective is to recover the surface normals \mathbf{N} and depth z .

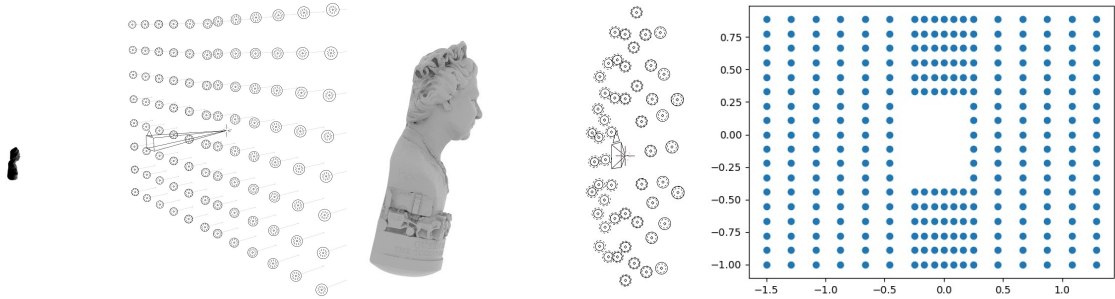


Figure 5: On the left the disposition of point-lights for the DiLiGenT dataset [52]. On the middle the one for the LUCES dataset [37]. In order to give idea of their scale, both configurations have been rendered facing the same object Queen at exact distance/orientation from the camera/light setup used for generating the PS images in the respective datasets (real LUCES and synthetic DiLiGenT). On the right, a potential (out of the ones sampled at train time) point light configuration is shown. This is a rectangle of sides 3x2 with a hole in the middle of size 0.33x0.66 (sizing is in terms of z) containing the maximum of 288 lights.

3.2 Normal Prediction

The first step of our method includes training a CNN for per-pixel normal prediction using BRDF samples. This is done through the the observational map parameterisation introduced by [21] in order to tackle the far-field photometric stereo problem. Note that this is equivalent to BRDF inversion under the special case of $\hat{\mathbf{V}} = [0, 0, 1]$. As described in [21] an observational map records relative pixel intensities (BRDF samples) on a 2D grid (e.g. 32×32) of discretised light directions. Such a representation is highly convenient for use with classical CNN architectures as it provides a 2D input and is of fixed shape despite a potentially varying number of lights. While [21] proposes to train CNNs on rendered images of objects, it is shown in [31] that simpler per-pixel renderers can be used instead, making the training procedure much faster and simpler. We use the latter approach in this work. Following [31], an RGB observation map O_{rgb} of size $d \times d \times 3$ is constructed as:

$$O_{\text{rgb}} \left(\left\lfloor \frac{d \hat{L}_m^x + 1}{2} \right\rfloor, \left\lfloor \frac{d \hat{L}_m^y + 1}{2} \right\rfloor \right) = \begin{bmatrix} j_r / \phi_r \\ j_g / \phi_g \\ j_b / \phi_b \end{bmatrix}_m \quad (4)$$

In addition, we note that in the case of specular reflection, the BRDF samples j are dependent on the viewing vector \mathbf{V} . This variation is only expected to be significant in the case of perspective projection for points not close to the imaging center. Nonetheless, the set of orthographically rendered observation maps considered in [21] or [31] is only a special case of the possible observation maps. Thus, to make the network training problem easier, we extend the observation map concept to incorporated the viewing vector \mathbf{V} (which is known and constant for all light sources m) such as:

$$O = [O_{\text{rgb}} ; \mathbf{1V}] \quad (5)$$

where $\mathbf{1}$ is $d \times d \times 3$ and $;$ is a concatenation on the 3rd axis so defining a $d \times d \times 6$ map. Finally, these observation maps are fed into a CNN which regresses surface normal \mathbf{N}_p . The CNN is trained with the angular loss defined as $|\text{atan2}(\|\mathbf{N}_t \times \mathbf{N}_p\|, \mathbf{N}_t \cdot \mathbf{N}_p)|$ with \mathbf{N}_p are the predicted normals and \mathbf{N}_t are the ground truth normals.

3.3 Adapting to the Point Light Setup

In order to solve the point light PS problem for a realistic capture setup we adapt the training procedure to only sample observation maps which are plausible at test time. Therefore, instead of sampling a random set of light directions as in [31], we sample 3D points inside a virtual camera frustum. For each point, a different LED configuration is simulated.

Configuration sampling. The sampling procedure for a point begins with sampling normalised image plane coordinates $[u, v] \in [-1, 1]$. Then camera focal length $f \in [1, 10]$ is sampled. $f = 1$ corresponds

to a real fish eye lens and $f = 10$ is close to orthographic viewing. For reference LUCES normalised focal length is ≈ 1.5 , DiLiGenT is ≈ 5 . Then a depth z is sampled in a range from 10cm to 170cm and this depth is used to back-project image coordinates and obtain 3D point in camera coordinate system $\mathbf{X} = [uz/f, vz/f, z]^T$.

The rest of the configuration is sampled proportionally to z which allows⁴ for tackling LED arrangement of vastly different scales (see Figure 5). We assumed that point lights are approximately on a plane parallel to XY axes. The plane offset is uniformly sampled in the range $[0, 0.25z]$ and all lights are positioned at a height with respect to that plane up to $\pm 0.05z$. In terms of distribution of the lights on that plane, we assume a rectangle with side lengths in $[0.5z, 3z]$ and with a rectangular hole in the middle with side lengths $[0, 0.66z]$ (see Figure 5). This plane area is divided into a grid and a number between 15 and 288 points of this grid are selected to be the light positions \mathbf{P}_m . Light brightness ϕ_m , are sampled uniformly and independently in log scale from $\phi_m \in [0.25, 4]$, dissipation factors $\mu \in [0, 3]$ and $\mathbf{D}_m = [d_x, d_y, 1 + d_z]$ with $d_{x,y,z} \in [-0.1, 0.1]$ (ensuring $\|\mathbf{D}_m\| = 1$). Finally, surface normal \mathbf{N} , material parameters and global illumination approximations are sampled independently following the exact same hyper-parameters of [31].

Reflectance rendering. Once the point parameters are sampled, the image samples i_m are rendered (Equation 2) using the renderer of [31] with the addition of point light attenuation (Equation 1). Additional global illumination approximations for shadows, reflections and ambient light are also applied as in [31]. Schematically this corresponds to:

$$\{\mathbf{X}, \mathbf{P}_m, \phi_m, \mathbf{D}_m, \mu_m, \mathbf{N}\} \xrightarrow{\text{Eq.1}} \{\hat{\mathbf{L}}_m, a_m\} \xrightarrow{\mathbf{N}, \text{Render}} \{i_m\}. \quad (6)$$

Note rendered intensities i_m are rendered using constant light attenuation $a_m = \phi_m$. The final intensities i_m are obtained by using non-linear radial model of dissipation described in Eq. 1. Also note that 10bit discretisation, and saturation are applied (i.e. conversion to integer $\in \{0, 1023\}$ and re-normalisation) when rendering values $\{i_m\}$ to approximate the camera of LUCES.

Observation map generation. After performing the point rendering to compute i_m , the aim is to compensate for light attenuation to compute reflectance sample (Equation 3) and generate observation maps (Equation 4). In order to get robustness to imprecise depth initialisation at test time, the training procedure involves perturbing the ground truth depth value z by $\delta z \sim \mathcal{N}(0, 5\%z)$ ⁵ to obtain $z' = z + \delta z$. In addition, all setup parameters are also slightly perturbed to account for potential setup miss-calibration i.e.:

$$\{z, \mathbf{X}, \mathbf{P}_m, \phi_m, \mathbf{D}_m, \mu_m\} \xrightarrow{\delta} \{z', \mathbf{X}', \mathbf{P}'_m, \phi'_m, \mathbf{D}'_m, \mu'_m\} \quad (7)$$

The hyper-parameters in Equation 7 are:

$\delta \mathbf{P} \in [-0.1\%z, 0.1\%z]$ (additive), $\delta \phi \in [0, 1\%]$ (multiplicative), $\delta \mathbf{D} \in [-0.1, 0.1]$ (additive), $\delta \mu_1 \in [0, 0.1]$ (additive) and $\delta \mu_2 \in [0, 10\%]$ (multiplicative). We note that we samples these perturbations both independently for all light sources but also include an additional amount of perturbation (of same distribution) to all the lights at the same time to account for systematic errors. Finally, these perturbed parameters are used to recompute attenuation, reflectance samples and then observation maps O' i.e.:

$$\{\mathbf{X}', \mathbf{P}'_m, \phi'_m, \mathbf{D}'_m, \mu'_m\} \xrightarrow{\text{Eq.1}} \{\hat{\mathbf{L}}'_m, a'_m\} \xrightarrow{i_m, \text{Eq.3-5}} O' \quad (8)$$

Specific setup. In the case of aiming to train a network for a specific dataset (e.g LUCES [37]) the plausible observation map space is highly reduced since the camera/light configuration is known, and thus the data generation process can take advantage of that. We note that setup knowledge means that the parameters used to compute the observation maps at test time, i.e $\{\mathbf{P}'_m, \phi'_m, \mathbf{D}'_m, \mu'_m\}$ (in Equation 8) are assumed to be known at train time. Of course, it is still desirable to have some robustness to potential setup miss-calibration, therefore the perturbation equation is applied in *reverse* order (so as to end up at the map creation stage with the parameters that will be used at test time). This is summarised as:

⁴We assume that in most cases the radius of the arrangements of point lights would be of similar scale as the width of the objects which is being scanned.

⁵The Gaussian distribution encourages the network to be more accurate when δz is small and thus get an improvement in the iterative setting.

$$\text{Calibration} \xrightarrow{\text{Copy}} \{\mathbf{P}'_m, \phi'_m, \mathbf{D}'_m \mu'_m\} \quad (9)$$

$$\{\mathbf{P}'_m, \phi'_m, \mathbf{D}'_m \mu'_m\} \xrightarrow{\delta} \{\mathbf{P}_m, \phi_m, \mathbf{D}_m \mu_m\} \quad (10)$$

$$\text{Eq. 6, 8} \rightarrow O' \quad (11)$$

3.4 Iterative Refinement of Depth and Normals

Assuming an estimate of normals, the depth can be obtained by numerical integration. This is performed using the ℓ_1 method of [45]. The variational optimisation includes a Tikhonov regulariser with $z = z_0$ (z_0 is the previous estimate of the depth map starting from plane) weight $\lambda = 10^{-6}$ and is solved in a ADMM scheme⁶.

As the BRDF samples j (see Equation 3) depend on the unknown depth, they cannot be directly computed to be input to the network. To overcome this issue, we employ an iterative scheme where the previous estimate of the geometry is used. The procedure involves computing the near to far conversion as described, obtaining a new normal map estimate through the CNN and finally numerical integration. See Figure 4 for an example of intermediate results of our iterative procedure. As it is the case in competing classical methods [32, 44], this iterative procedure is initialised with a flat plane at the approximate mean distance.

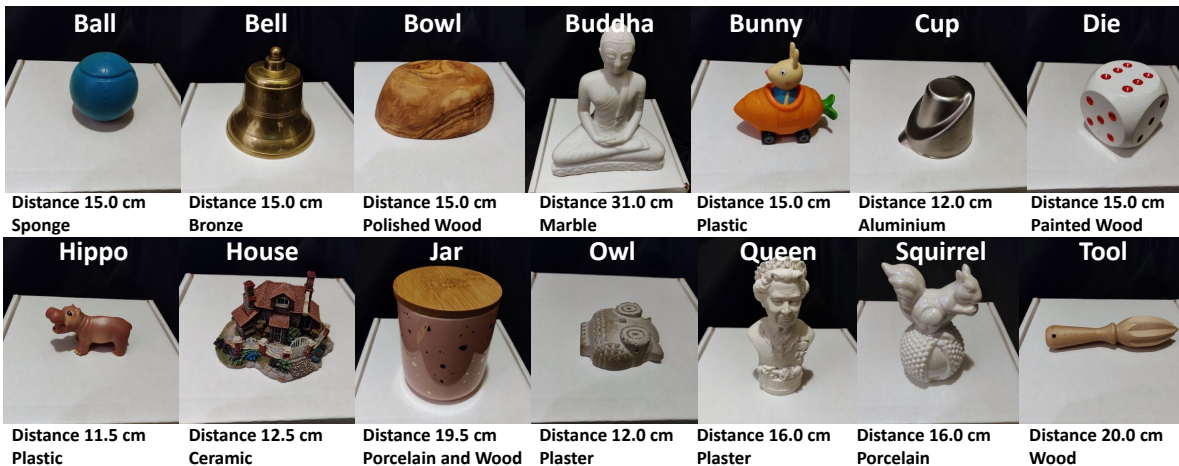


Figure 6: Top view of the objects captured for LUCES dataset. Below every object the acquisition distance between the object and the camera, and the material of the object are reported.

4 LUCES Dataset

This section gives an overview of the data capture and calibration procedure for the LUCES dataset, first presented in [37].

The Photometric stereo setup. Our setup (see Figure 2, left) consists of the following main components:

- RGB camera FLIR bfs-u3-32s4c-c with 8mm lens,
- 52 LED Golden Dragon OSRAM,
- variable voltage for adjustable LED power,
- Arduino Mega 2560.

A custom printed circuit board (PCB) has been designed to host 52 bright LED controlled with by an Arduino Mega. The configuration of the LEDs was planar around the camera. A set of 52

⁶Code ported from https://github.com/yqueau/normal_integration.



Figure 7: Demonstration of the processing steps performed per object in LUCES dataset. Firstly, compensation for radial distortion and demosaicing is performed on raw images to get RGB ones (left). CT-scanned ground truth meshes are aligned with RGB images and ground truth normal maps are rendered (middle). Segmentation masks are also manually generated.

images was captured per object. The camera parameters (aperture and shutter speed) and LED voltage were adjusted to achieve the best object exposure, which is very critical for specular objects. All camera preprocessing was turned off during the acquisition, including white-balance and analog gain. Several optomechanical tools have been used for holding the camera and the PCB jointly. A manual XYZ translation stage with differential adjusters has been used to positioning the camera accurately through the printed circuit board. In order to limit interreflections and ambient light, the walls surrounding the setup have been covered with black, polyurethane-coated nylon fabric.

Camera intrinsics. This is performed using 100 checkerboard images and the OpenCV calibration toolbox. Fourth degree radial distortion is estimated and this is used to rectify all the images. The calibration re-projection error was 0.42px. The RAW data (before demosaicing and rectification) are also available.

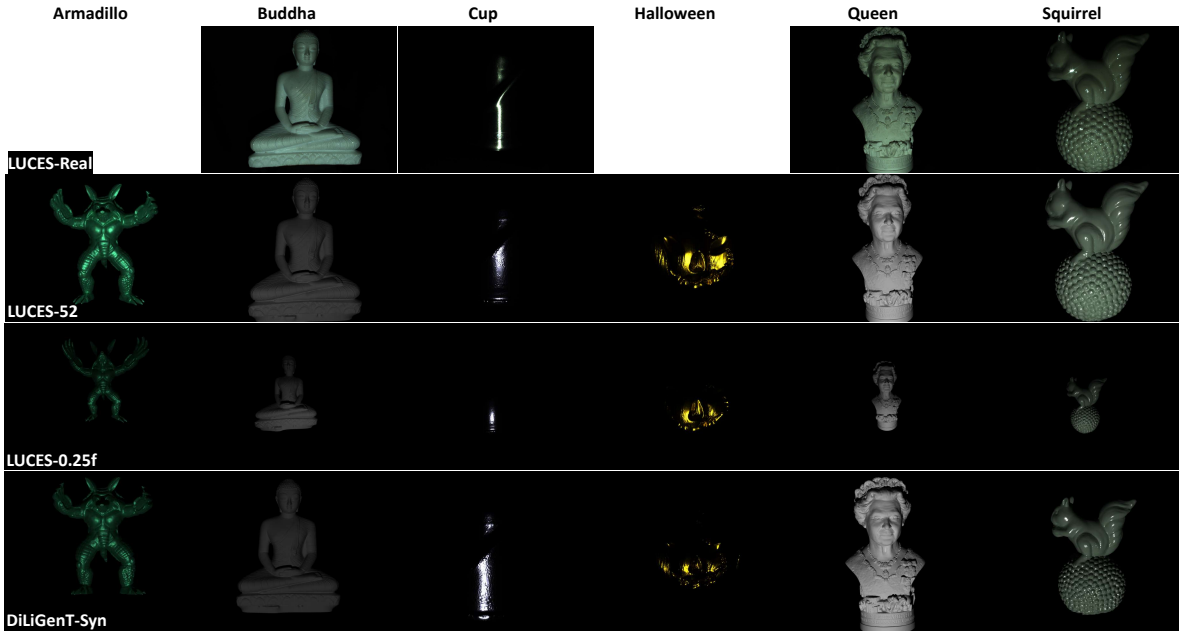


Figure 8: One image example for all synthetic objects rendered as well as the corresponding real ones in top row (except for Armadillo and Halloween that have no real counterparts). The second row aims to closely match the configuration of LUCES [37], the third row reduces the focal length by a factor of 4 and the bottom row aims to closely match DiLiGenT [52].

Light calibration. This section presents the method used to estimate all the point light parameters introduced in Section 3.1 ($\{\mathbf{P}_m, \phi_m, \mathbf{D}_m, \mu_m\}$). To do so, we captured PS images of a purely diffuse reflectance plane i.e. 99% nominal reflectance in UV-VIS-NIR wavelength range (350 - 1600nm). To have an initial estimate of ϕ_m , we measured the brightness of the LEDs with a LuxMeter. For every object, the calibration plane was captured twice, at different distances, in order to get data redundancy and produce a more accurate calibration. Thus, the Lambertian calibration object with albedo ρ and

surface normal \mathbf{N} , should satisfy the resulting image irradiance equation:

$$i_m = \phi_m a_m \rho \hat{\mathbf{L}}_m \cdot \hat{\mathbf{N}}. \tag{12}$$

The irradiance Equation 12 was implemented into a differentiable renderer (using Keras of Tensorflow v2.0) with the LED parameters being the model weights thus allowing refinement from a reasonable initial estimate. The parameters were initialised as follows: ϕ_m from the LuxMeter, $\mathbf{D}_m = [0, 0, 1]$, $\mu_m = 0.5$, \mathbf{P}_m from the schematic of the printed circuit board of the LEDs and $\rho = 0.5$. We used L_1 loss function for 30 epochs and converged to around 0.005 error i.e 0.5% of the maximum image intensity. The complete calibration parameters are included in the dataset.

3D ground truth capture. Initial version of the 3D ground-truth [37] was acquired with the optical 3D scanner GOM ATOS Core 80/135 with a reported accuracy of 0.03mm (see Figure 2). The GOM scanner uses a stereo camera set-up and more than a dozen scans were performed and fused per object. In order to keep the geometry of the object consistent with the PS data, no spray coating has been used to ease the acquisition. In this work, we provide more accurate ground truth meshes for all the non-diffuse objects⁷.

Alignment. The scans of the objects were aligned and merged using MeshLab [12]. Some manual removal of noisy regions was performed and finally Poisson reconstruction was used in order to obtain full continuous surfaces which are both useful for rendering normal maps and for mutual information alignment. As expected, not all parts of the surfaces of all objects have the same amount of noise, especially the metallic objects (Bell, Cup). Meshes were aligned with the photometric stereo images using the mutual information registration filter of MeshLab [12]. This was initialised manually and pixel perfect accuracy was achieved. Using the aligned meshes, ground truth normal maps were rendered (using Blender). In addition, manual segmentation was performed to remove regions where the GT was unreliable (markers on the objects, holes etc). The steps per object are summarised in Figure 7.

Dataset overview. For each of the 14 objects (see Figure 6), 52 PS images have been acquired using the BayerRG16 RAW format. The total amount of PS images amounts to 728. For all objects, rectified RGB PS images are available. We note that color balancing was not performed on the images as this distorts the saturated pixels and ultimately loses information. Instead, RGB light source brightness are provided along with the rest of point light source parameters. Both normal map and depth ground truth are provided in order to evaluate the accuracy of near-field PS methods with either case.

Non GT objects. We also captured 15 light images of 3 additional objects shown in Figure 13. Metallic-silver Bulldog statue, a porcelain Frog as well as a mutli-object scene featuring a shiny wooden elephant statue in front of the porcelain Squirrel. Bulldog and Frog were too hard to scan and the elephant and squirrel could not be transported in their exact configuration to the CT scanner.

5 Experimental Setup

In this section we provide various experimental setup details related to CNN training and datasets used for evaluation.

5.1 CNN Training

We use the exact architecture of PX-NET [31] which is a miniature version of DenseNet [19] with 4.5M parameters. We trained 3 networks, one with general data data and 2 specific ones, one for the LUCES configurations and one for the DiLiGenT one. For the general one, we trained for 50 epochs and selected the checkpoint with best DiLiGenT performance (35). The MAE evolution is shown in Figure 9 for this network. The specific networks converged more quickly taking 9 epochs for LUCES and 25 epochs for DiLiGenT. A batch size of 2400 and 5000 steps per epoch was used in all experiments. It took 20 minutes to complete an epoch on a machine with 2 Titan RTX GPUs.

⁷Scans of all the objects except Buddha, House, Owl, Queen were obtained with the Zeiss CT scanner M1500/225 kV which provides an accuracy within the order of $9\mu m$.

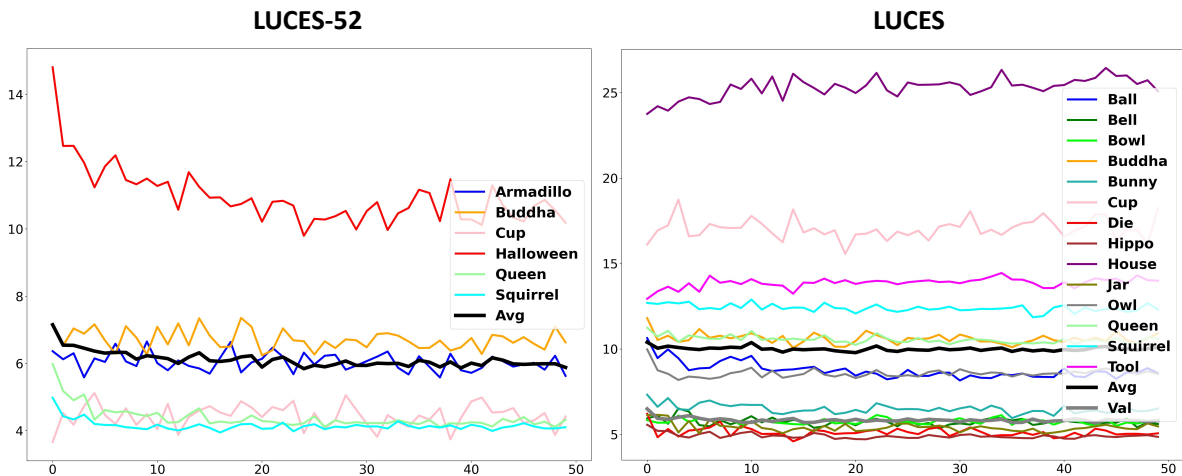


Figure 9: MAE evolution (during training) curves illustrate the performance of our setup-agnostic network on predicting normals (NfCNN) for synthetic (left) and real data (right). We note that we used the real DiLiGenT dataset as validation loss and selected the checkpoint (35) where this is minimised. We observe that although the average error is gradually decreasing, for some real objects (House, Cup) the performance is actually getting worse with more training signifying that some real effect is not properly modelled.

5.2 Datasets

We evaluate our method on the real datasets LUCES (see Section 4) and DiLiGenT [52]. Additional evaluation is performed on four synthetic datasets namely LUCES-52, LUCES-15, LUCES-0.25f and synthetic DiLiGenT. More details about each dataset are provided below:

DiLiGenT [52]. It contains 96 images of the resolution of 612×512 px for each of 10 captured objects. It is usually assumed that this is a far field dataset with the directional uniform illumination. However, in reality LEDs were used for the illumination, and their positions \mathbf{P}_m , brightness ϕ_m and perspective camera parameters are also provided. Using the light positions \mathbf{P}_m , for each object the mean distance can be computed to match the average light directions $\hat{\mathbf{L}}_m$. Finally, LED directions \mathbf{D}_m were assumed perpendicular to the LED plane, and $\mu = 0.5$ was also assumed.

Synthetic LUCES-52 [37]. In order to have an estimate of the synthetic to real gap, we chose 4 objects from LUCES and rendered them in exactly the same position and similar materials: Buddha, Queen - Lambertian, Cup - metallic, Squirrel - specular dielectric. We also rendered two additional synthetic objects: Armadillo and Halloween. Armadillo was chosen as it has a challenging geometry (with occlusion boundaries at hands/face) and was rendered with an ‘intermediate’ Disney [8] material (all parameters⁸ set to 0.5). The Halloween object was chosen to be rendered metallic-gold as it is very hard to laser scan real metallic objects with concavities and other high frequency surface details. All objects were rendered with the Cycles rendering engine of Blender [6] to generate realistic global shadows and self reflections. The default path tracing integrator of Cycles was used using 100 samples per pixel, 8 light bounces as well as no post-processing de-noising.

Synthetic LUCES-15 and LUCES-0.25f. In addition to the synthetic version of LUCES described above, we consider another two variations namely Synthetic LUCES-15 and Synthetic LUCES-0.25f. The first one simply contained 15 out of the 52 images per object and was aimed at providing an evaluation at a sparse lighting setting which is usually preferred in practice. LUCES-0.25f was rendered with exactly the same objects and materials by reducing the camera focal length from 8mm to 2mm, in order to simulate a fish-eye lens. Object sizes/positions were also adjusted to keep them aligned in the middle of the field of view. The purpose of this dataset is to test a situation where perspective viewing becomes significant. See Figure 8 for example images.

Synthetic DiLiGenT. Finally, the same 6 synthetic objects were rendered in the DiLiGenT [52] configuration. The aim of this experiment was to assess potential performance improvement when the number of lights increases from 52 to 96 and the capture setup becomes for ‘far-field’ - higher distance

⁸Due to the highly non-linear nature of the BRDF, this material is not necessarily the average brightness.

from the camera and higher focal length.

5.3 Evaluation Protocol

This section describes the evaluation protocol for all the above mentioned datasets.

Competitors. We compare our method against the far-field CNN approaches of [21], [31], the near-field variational optimisation methods [32], [44] as well as the recent near-field CNN-based method of [50]. For all 3 CNN-based methods, the same network checkpoint was used as the one in the corresponding papers. For all methods, test code was available online⁹. Finally, for the comparison on the DiLiGenT [52] benchmark (see Table 3), we also report the numbers of some other competing approaches.

Naive vs Compensated usage of far-field methods. In order to demonstrate the importance of our point-light compensation procedure, we compare the usage of the far-field CNN approaches of [21], [31] with/without using it. Naive usage refers to using the raw image values (i.e. i_m) without attenuation compensation for computing observation maps as well as the average light direction for each LED. The predicted normals are also integrated using our method in order to have a qualitative shape comparison.

Evaluation metrics. For most experiments, the evaluation metrics are mean angular error (MAE) on normals (in degrees) as well as mean depth error (in mm) on computed depth maps. We note that the real DiLiGenT [52] benchmark only reports ground truth normals and also for 3 of our objects, no ground truth was available so the comparison is only qualitative (see Figure 13). In addition, we note that the variational optimisation methods [32], [44] only output depth maps, therefore in order to have comparison in the normal domain for them, normal maps are generated with numerical differentiation. Therefore, for the rest of the methods we report 2 types of normal maps namely NfCNN (normals from CNN-network predictions) and NfS (normals from shape-obtained though numerical differentiation).

Input resolution. Most LUCES experiments (both real and synthetic), were performed at a quarter resolution i.e. 512×384 px in order to have a fair comparison with [50] which is GPU memory limited (even on their original paper the authors report unable to run on more than 600×600 px resolution on a 48GB GPU RAM). However, it has to be noted that we are unsure if some of their respective hyperparameters is resolution dependent. For the real LUCES data, we also present our evaluation on full resolution (2048×1536 px) images and show that our method is resolution independent (with around $0.1mm$, 0.1° difference between quarter and full scale). DiLiGenT (both synthetic and real) on the other hand offers maximum resolution of 612×512 px which was used for all the relevant experiments reported.

6 Experiments

In this section we present the results of the various experiments on the synthetic and real datasets introduced in the previous section.

6.1 Synthetic Experiments

This section explains the synthetic experiments which are summarised in Table 1. A further breakdown per category follows.

Shape integration. The first experiment we conducted aimed at calibrating the quality of the numerical integration of the normal map. As no realistic depth map is C2 continuous, GT normals are not compatible with the GT depth. Indeed, integrating the GT normals and then re-calculating them with numerical differentiation introduces 3.52° MAE on average (Table 1 top) on full resolution and even reaches 6.50° at half resolution as the naive numerical differentiation is very resolution dependent.

⁹ [32] https://github.com/fotlogo/semi_calib_ps_cvpr2017

[44] https://github.com/yqueau/near_ps

[50] <https://github.com/hiroaki-santo/deep-near-light-photometric-stereo>

[21] <https://github.com/satoshi-ikehata/CNN-PS>

Method	LUCES-52						LUCES-0.25f		LUCES-15		DiLiGenT-Syn									
	Armadillo	Buddha	Cup	Halloween	Queen	Squirrel	AVG Norm	AVG Depth	AVG Norm	AVG Depth	AVG Norm	AVG Depth	Armadillo	Buddha	Cup	Halloween	Queen	Squirrel	AVG Norm	AVG Depth
GT Normals - NfS	4.42	2.69	2.85	3.7	4.22	3.26	3.52	3.25	4.14	2.17	4.14	2.17	7.77	6.34	4.35	6.55	9.40	8.60	7.17	3.93
GT Normals (half resolution) - NfS	7.47	6.75	1.65	6.16	10.12	6.88	6.50	3.36	9.18	2.29	9.18	2.29	-	-	-	-	-	-	-	-
Naive CNN-PS - NfCNN	18.11	22.58	7.83	20.54	17.24	16.25	17.09	-	22.59	-	19.82	-	5.87	8.43	8.50	25.50	7.46	6.80	10.43	-
Naive - CNN-PS - NfS	18.81	22.94	7.26	20.48	18.24	16.77	17.42	6.53	22.05	8.23	19.78	7.92	7.72	9.02	7.16	22.01	8.48	8.04	10.40	4.93
Naive PX-NET - NfCNN	15.7	23.01	10.10	12.96	19.30	17.82	16.48	-	22.6	-	17.66	-	4.91	8.37	3.88	10.53	6.59	5.74	6.67	-
Naive - PX-NET - NfS	16.13	23.37	9.85	12.85	20.38	18.24	16.80	7.22	22.37	9.09	17.87	7.57	7.33	8.88	4.28	10.93	8.03	7.37	7.80	4.37
Compensated CNN-PS - NfCNN	16.18	33.30	6.93	22.55	8.06	7.72	15.79	-	11.66	-	19.28	-	5.67	7.70	8.61	21.82	5.99	5.73	9.25	-
Compensated CNN-PS - NfS	16.39	27.88	6.69	22.12	9.92	9.24	15.38	6.05	12.54	3.87	17.6	7.71	7.49	8.30	7.42	19.38	7.24	7.23	9.51	4.68
Compensated PX-NET - NfCNN	6.64	6.14	4.93	11.16	4.7	4.48	6.34	-	7.6	-	7.56	-	4.06	5.66	2.90	10.58	4.79	4.15	5.36	-
Compensated PX-NET - NfS	8.48	7.48	4.97	11.91	7.87	6.94	7.94	2.65	9.73	2.5	2.76	8.85	6.78	6.46	3.37	11.03	6.70	6.27	6.77	3.61
Ours - GT Depth - NfCNN	6.21	6.41	4.45	11.13	4.21	4.08	6.08	-	5.61	-	7.08	-	3.43	6.63	3.49	11.81	4.39	4.38	5.69	-
Ours - Iteration 1 - NfCNN	6.21	6.47	4.54	11.56	4.24	4.11	6.19	-	6.03	-	7.21	-	3.42	6.64	3.49	11.95	4.39	4.38	5.71	-
Ours - Iteration 1 - NfS	8.51	7.79	4.74	12.46	7.68	6.65	7.97	3.04	8.45	2.49	8.59	3.05	6.38	7.15	4.10	12.86	6.53	6.40	7.24	4.3
Ours - Iteration 2 - NfCNN	6.22	6.40	4.45	11.24	4.22	4.08	6.10	-	5.66	-	7.1	-	3.42	6.63	3.49	11.90	4.40	4.38	5.70	-
Ours - Iteration 2 - NfS	8.5	7.69	4.71	12.17	7.68	6.65	7.90	2.92	8.17	2.31	8.49	2.9	6.38	7.14	4.10	12.81	6.53	6.40	7.23	4.3
Quéau et. al	22.45	12.95	22.02	38.19	12.11	11.02	19.79	7.87	17.47	5.99	20.17	8.32	16.16	8.65	13.33	30.81	10.99	11.97	15.32	9.15
Logothetis et. al	23.32	11.59	35.48	49.61	13.66	13.46	24.52	5.03	27.79	5.58	24.9	5.21	20.79	8.16	38.86	41.75	11.78	15.03	22.73	6.28
Santo et. al	15.79	17.52	5.65	21.41	12.57	10.70	13.94	4.38	13.36	4.17	15.52	4.46	7.52	11.63	8.97	20.14	8.28	6.56	10.52	6.8

Table 1: Full quantitative comparison on synthetic data. For our method, we report raw normal prediction error as NfCNN and numerically differentiated normals as NfS. Also report NfCNN error when the GT depth is used as initialisation and also NfS error for integrating the GT normals. We compare Logothetis et. al [32] and Quéau et. al [44], Santo et. al [50], CNN-PS [21] and PX-NET [31]. For LUCES-52 and synthetic DiLiGenT, metrics on all objects are shown; for the other 2 only average errors are reported (the aim is to observe drop of performance by reducing the number of lights or focal length). Note that for the NfCNN-GT reported metric, only the normal error is meaningful.

Therefore, we do not compare raw network predictions (NfCNN) and MAE after differentiation of the surface (NfS) and expect the first figure to be lower for networks trained to regress normals¹⁰.

Naive usage of far-field networks. The next experiment consists of naively using the far-field methods [21, 31] with no point light compensation. As expected, in all ‘near field’ results (except in synthetic DiLiGenT), both normal and depth errors are significantly higher than all other competitors and this is better understood visually in Figure 10; the shape is ‘locally correct’ with no bumps at specular highlights or other similar artifacts but still severely distorted. To add to this point, using these networks as part of our iterative process (marked as Compensated) drastically improves the result and in fact they outperform the variational optimisation methods of [32] and [44]. The difference between naive and compensated is significantly lower on the synthetic DiLiGenT as scale of this dataset is mostly far-field with the point light effect being less important.

Proposed method. We report results of our method using the setup agnostic network. We report error at iteration 1 (i.e. compensation using the initial planar geometry estimate) and iteration 2 and observe a marginal improvement signifying that the process has converged. We also report NfCNN where the point light compensation has been performed with the GT depth to estimate the limiting performance of the iterative method. We again confirm that this is only marginally better than iteration 2 error confirming that our network is not very sensitive to depth initialisation.

Material variation. We note that the proposed method performs similarly in all 6 objects despite the significant material variations. This is not the case for the variational optimisation competitors ([32], [44]) which are significantly worse on the metallic objects (Cup, Halloween) than the Lambertian ones (Buddha, Queen). Quite surprisingly, [50] performs worse on the Lambertian objects possibly signifying lack of training data with exact Lambertian reflections.

LUCES-0.25f. We observe no drop of performance between LUCES-52 and its lower focal length counterpart verifying that our method correctly compensates for the effect of perspective viewing. In contrast, for the naive far-field methods the error is increased.

LUCES-15. We observe small drop of performance in the normal error between LUCES-52 and the 15 lights version (7.1° to 6.1°) but the depth error remains practically the same - $2.9mm$. This is probably explained by the fact that in the low light setting a few points become unsolvable (inflating the mean

¹⁰Some more sophisticated shape differentiation method such as [68] would probably reduce the discrepancy between these 2 metrics but that would not alter the discussion of this section.

Method	Error	Ball	Bell	Bowl	Buddha	Bunny	Cup	Die	Hippo	House	Jar	Owl	Queen	Squirrel	Tool	Average
Logothetis et. al [32]	MAE	12.55	36.02	16.66	14.59	13.01	28.08	11.54	14.27	32.20	12.21	16.61	15.99	17.73	15.72	18.37
	MZE	1.60	5.60	6.11	5.49	2.14	3.17	4.35	2.09	7.62	6.08	4.06	5.11	2.10	5.34	4.35
Quéau et. al [44]	MAE	8.96	28.16	16.16	15.67	12.60	43.20	9.49	13.93	33.01	10.30	17.39	16.01	16.55	16.86	18.45
	MZE	3.67	12.83	6.54	7.72	2.56	13.03	4.14	8.21	6.19	11.97	5.51	5.92	4.67	6.93	7.14
Santo et. al [50]	MAE	13.27	10.03	7.27	19.22	9.44	11.74	4.95	7.20	31.49	10.00	12.09	13.43	13.24	10.90	12.45
	MZE	2.48	2.55	3.94	6.24	5.13	1.44	2.65	3.76	8.13	7.11	5.06	5.47	2.07	5.96	4.43
Naive PX-NET [31]	MAE	14.23	13.47	8.54	29.99	18.99	13.84	15.59	21.17	35.23	19.38	21.19	21.77	25.63	19.33	19.88
	MZE	4.16	4.37	7.20	18.22	3.36	1.36	5.25	5.46	8.74	11.15	2.88	5.56	3.76	3.83	6.09
US - Compensated PX-NET	MAE	10.90	8.99	6.03	12.27	7.77	14.46	8.02	11.92	29.80	7.00	17.69	11.25	13.87	13.44	12.39
	MZE	0.80	2.03	2.16	3.80	2.55	1.60	1.81	3.30	8.79	4.44	4.60	3.64	1.54	2.17	3.09
US - Specific	MAE	8.62	6.36	6.41	11.59	6.85	14.09	4.76	5.67	22.53	5.63	8.86	9.82	11.84	11.19	9.59
	MZE	0.59	1.66	2.22	3.54	2.39	1.84	1.55	1.51	8.34	2.32	4.23	3.12	1.15	5.12	2.83
US - General	MAE	8.86	7.51	5.96	11.60	7.07	15.27	5.19	5.54	22.91	6.14	8.86	9.96	11.77	11.56	9.87
	MZE	0.58	1.72	2.07	3.82	2.33	2.16	1.86	1.85	8.84	2.71	4.28	2.85	0.76	3.67	2.82
US - General Full res.	MAE	8.84	7.51	5.95	11.59	7.06	15.35	5.19	5.60	22.97	6.19	8.89	9.97	11.77	11.64	9.90
	MZE	0.58	1.77	1.96	3.74	2.32	2.17	1.89	1.83	8.86	2.63	4.23	2.81	0.76	3.69	2.80

Table 2: Evaluation on the LUCES [37] benchmark compering with Logothetis et. al [32], Quéau et. al [44], Santo et. al [50], and PX-NET [31]. Mean angular error (MAE in degrees) of predicted normals and mean depth errors (MZE in mm) are reported.

error) but the overall surface can still be recovered. This is not the case for [50] where both normal and depth errors increase. The variational optimisation competitors ([32] and [44]) also have minimal drops of performance in this low light setting. A surprising result is that compensated PX-NET is also performing similarly between the 52 and 15 lights settings even though it was trained with a minimum of 50 lights. This is probably explained by the fact that it was trained to be very resilient to shadows which essentially reduce the amount of active lights.

6.2 Real Data evaluation.

This section presents the results of the real data evaluation on the LUCES [37] and DiLiGenT [52] benchmarks.

LUCES [37]. Table 2 shows the quantitative evaluation on LUCES with qualitative comparison through normal maps in Figure 11 and shapes in Figures 12 and 13. We achieve the best performance in all objects with the exception of the metallic Cup where [50] is the best performer. This may be due to the use of a patch-based network which is able to extract the more information from noisy metallic data. Finally, on the qualitative only data of Figure 13, we note that optimisation competitors ([32] and [44]) struggle at the metallic Bulldog. This is not the case for [50] which seems to struggle at the high curvature region of the Frog neck. The proposed method performs reasonably on all 3 objects.

Synthetic to real gap. We observe that we perform significantly worse on the real LUCES objects with respect to their synthetic counterparts. As the geometry and lights are similarly matched we conclude that more research is needed in modeling and sampling realistic materials as well as other potential corruptions of real images (better noise models). This is most evident for the metallic Cup where the normal error increases from 4.5° to 14.1° . However, for all CNN-based methods (ours, [31, 50]) the material’s specularly does not seem to be a significant factor of performance. Indeed, convex regions (where self reflections are negligible) are consistently recovered correctly regardless of the material: diffuse head of Queen, bronze Bell, plastic Hippo, wooden Bowl; with the only exception being the aluminium Cup. This is a clear advantage of CNN methods against the classical ones that require diffuse or mostly diffuse materials.

Normal vs depth errors. We notice that the normal predictions are more noisy as opposed to depth predictions. This could be due to noisy estimates of the normals from the ground truth meshes which is inevitable for any laser scanner (see in particular the Ball in Figure 11). As the ground truth depth is more reliable, it is a better evaluation metric compared to the ‘ground truth’ normals. See Figure 12 for depth evaluation.

Error distribution. We observe that the hardest regions are the ones containing high frequency details (sharp boundaries) such as House, bottom part of the Squirrel, details of the Queen, etc. An interesting observation is also that for [50] there is growing inaccuracy towards the external part of the reconstruction (see Bell, Cup and Jar in Figure 11) which is probably due to the orthographic camera assumption.

Method	Ball	Bear	Buddha	Cat	Cow	Goblet	Harvest	Pot1	Pot2	Reading	AVG
SPLINE-Net [67]	1.74	4.65	9.14	5.48	9.55	9.38	24.44	5.91	7.9	12.77	9.1
ICML [56]	1.47	5.79	10.36	5.44	6.32	11.47	22.59	6.09	7.76	11.03	8.83
Exemplars [20]	1.33	5.58	8.48	4.88	8.23	7.57	15.81	5.16	6.41	12.08	7.55
PS-FCN [17]	2.7	4.8	6.2	7.7	7.2	7.5	7.8	10.9	6.7	12.4	7.4
CNN-PS [21]	2.2	4.1	7.9	4.6	8	7.3	14	5.4	6	12.6	7.21
Inverse Model [59]	1.78	4.12	6.09	4.66	6.33	7.22	13.34	6.46	6.45	10.05	6.65
PX-NET [31]	2.0	3.6	7.6	4.4	4.7	6.9	13.1	5.1	5.1	10.3	6.28
Santo et. al [50]	5.8	5.87	9.92	6.44	7.75	8.66	15.68	7.84	9.86	13.87	9.17
Ours - Comp. PX-NET	1.83	3.04	6.95	3.26	4.86	5.96	13.77	4.03	4.78	10.00	5.85
Ours - General	2.62	2.87	6.85	3.52	5.98	5.97	12.24	4.20	5.24	9.39	5.89
Ours - Specific	1.60	2.96	7.34	3.32	4.74	5.72	12.42	4.18	5.00	9.34	5.66

Table 3: Quantitative comparison of the proposed method on the DiLiGenT benchmark [52]. The competitors are the far-field methods: SPLINE-Net [67], ICML [56], PS-FCN [17], CNN-PS [21], Inverse Model [59] and PX-NET [31] as well as the near-field method by Santo et. al [50].

DiLiGenT [52]. Final evaluation at Table 3. We note that even though this dataset is usually considered far-field with directional lights, our point light compensation procedure improves the performance of PX-NET [31] (the best performing far-field method) from 6.28° to 5.85° demonstrating the importance of point-light modeling in real data. It is also interesting and somewhat surprising that compensated PX-NET also outperforms our general network (5.85° vs 5.89°) and that signifies that perspective viewing (which is the most important difference of these networks) is not significant on this dataset as opposed to LUCES, as shown in Table 2. Finally, we note that the best performing method is the DiLiGenT-specific network which is not really surprising even though the margin is quite small (5.66° vs 5.85°).

Specific setup network. Finally we note that setup specific networks are marginally better than the setup agnostic one (which took more time to converge though) signifying that the light distribution is not a big challenge for the CNN.

7 Conclusion

In this work we presented a CNN-based approach tackling the point light Photometric Stereo problem in both near and far-field setting. We leveraged the capability of CNNs to learn to predict surface normals from reflectance samples for a wide variety of materials and under global illumination effects such as shadows and interreflection. Numerical integration is used to compute the depth from predicted normals and this in turn is used to compensate the input images to compute reflectance samples for the next iteration. Finally, in order to measure the performance of our approach for near-field point-light source PS data, we introduced the LUCES dataset containing 14 objects imaged in a configuration where attenuation due to point lights sources and perspective viewing are significant.

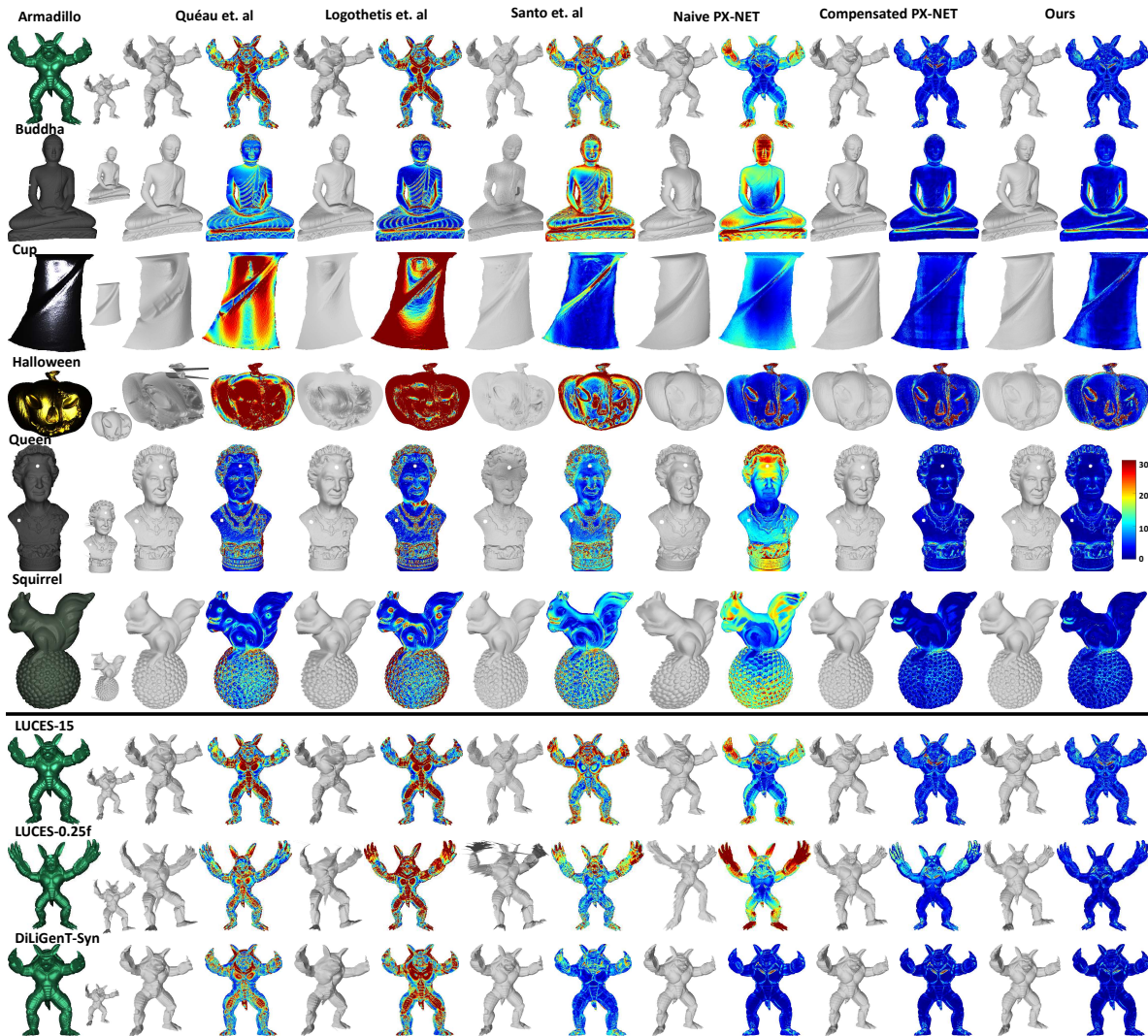


Figure 10: Visualisation of the results of Table 1 showing comparison with Logothetis et. al [32] and Quéau et. al [44], Santo et. al [50] and PX-NET [31] on all synthetic experiments. All 6 objects of synthetic LUCES-52 are shown and Armadillo is also shown for the other 3 synthetic datasets. For all objects, average PS image and GT depth shape is shown at the left. [32], [44] have very high error on the metallic objects (Cup, Halloween) as well as specular highlights on Squirrel and Armadillo. The naive far-field method PX-NET also has significant global deformation as it does not model the light attenuation effect. In contrast, the compensated version performs very well everywhere except the hands of Armadillo in the LUCES-0.25f due to its training data being rendered without perspective viewing. For [50], the error is concentrated towards the edges of each object as perspective projection is not modelled (this is especially evident on the LUCES-0.25f). The proposed approach achieves best performance in all objects.

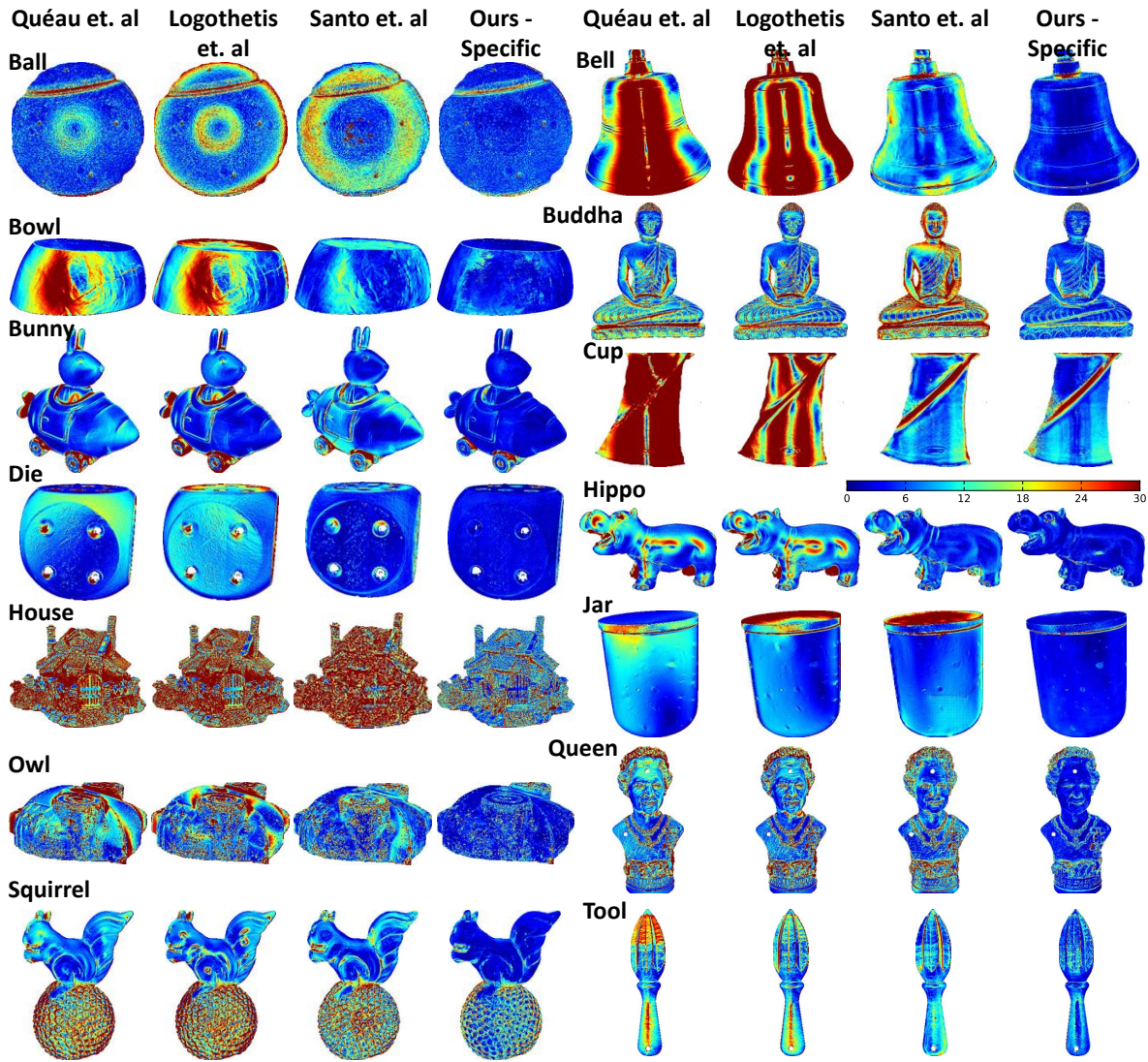


Figure 11: Visualisation of the normal error map comparison for all objects of real LUCES (Table 2) and all near-field methods: Logothetis et. al [32], Quéau et. al [44], Santo et. al [50].

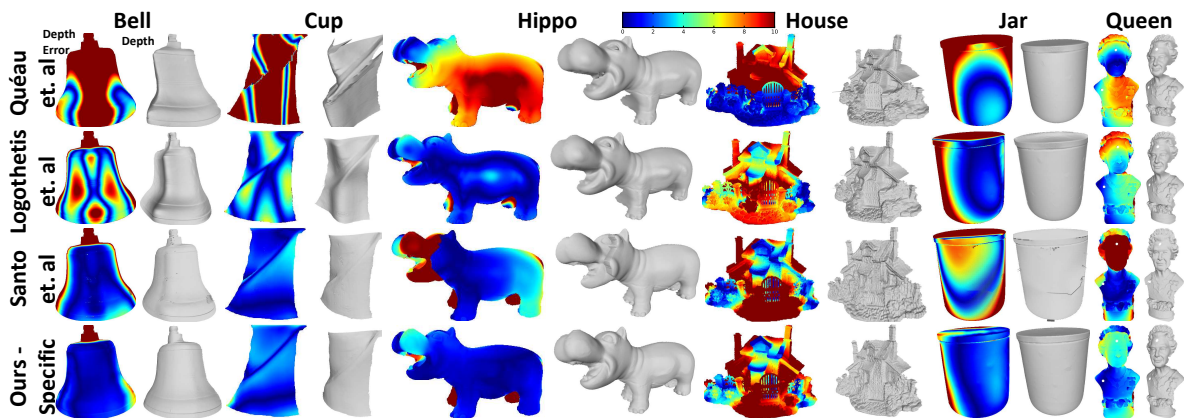


Figure 12: Output surface comparison for 6 objects of real LUCES (see Table 2 for quantitative results) and methods of Logothetis et. al [32], Quéau et. al [44], Santo et. al [50]. This is shown qualitatively through the 3D meshes and well as depth error maps (errors in mm).

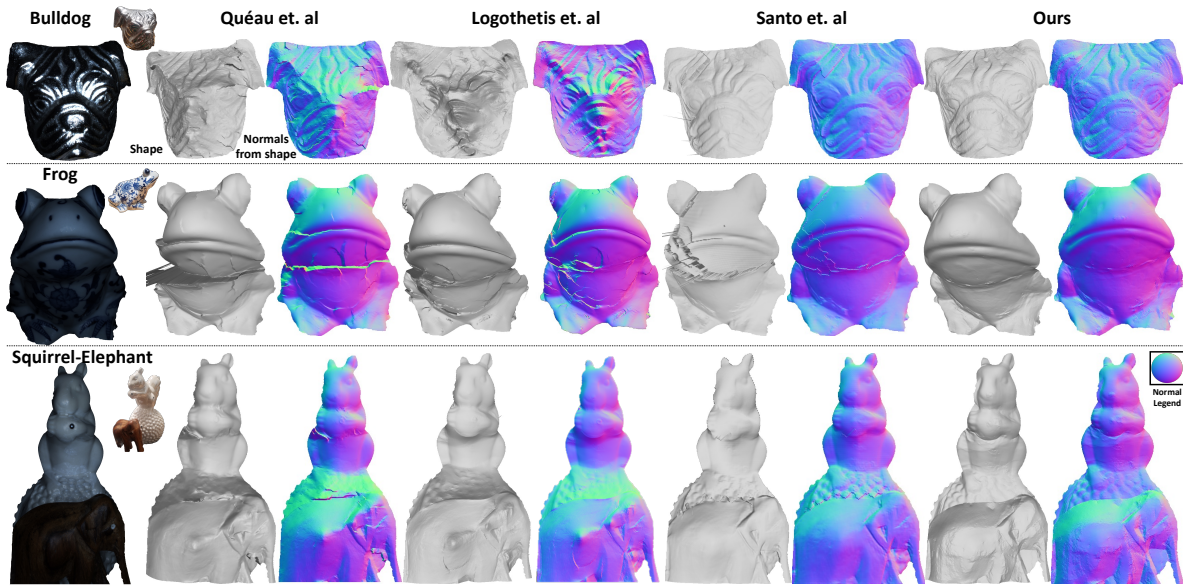


Figure 13: Qualitative comparison of the proposed method with [44], [32] and [50]. The first column shows the average Photometric Stereo image. In contrast to competition, the proposed approach has no visible deformation on the metallic object or the specular highlight in the middle of the elephant. In addition, there is a smooth recovery of belly of the Frog despite the shadows, as well as the bottom of Squirrel despite self reflection.

References

- [1] Aanæs, H., Dahl, A.L., Pedersen, K.S.: Interesting interest points - A comparative study of interest point performance on a unique data set. *International Journal of Computer Vision (IJCV)* **97**(1), 18–35 (2012)
- [2] Aanæs, H., Jensen, R.R., Vogiatzis, G., Tola, E., Dahl, A.B.: Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision (IJCV)* **120**(2), 153–168 (2016)
- [3] Ackermann, J., Goesele, M.: A survey of photometric stereo techniques. *Foundations and Trends in Computer Graphics and Vision* (2015)
- [4] Agrawal, A., Raskar, R., Chellappa, R.: What is the range of surface reconstructions from a gradient field? In: *European Conference on Computer Vision (ECCV)* (2006)
- [5] Alldrin, N., Zickler, T.E., Kriegman, D.J.: Photometric stereo with non-parametric and spatially-varying reflectance. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2008)
- [6] Blender-Online-Community: Blender - A 3D modelling and rendering package. Blender Foundation (2018). URL www.blender.org
- [7] Blinn, J.F.: Models of light reflection for computer synthesized pictures. In: *SIGGRAPH* (1977)
- [8] Burley, B.: Physically-based shading at disney. In: *SIGGRAPH Course Notes* (2012)
- [9] Chandraker, M.K., Agarwal, S., Kriegman, D.J.: Shadowcuts: Photometric stereo with shadows. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2007)
- [10] Chen, G., Han, K., Shi, B., Matsushita, Y., Wong, K.Y.K.: Self-calibrating deep photometric stereo networks. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2019)
- [11] Chen, G., Han, K., Wong, K.K.: PS-FCN: A flexible learning framework for photometric stereo. In: *European Conference on Computer Vision (ECCV)* (2018)
- [12] Cignoni, P., Callieri, M., Corsini, M., Dellepiane, M., Ganovelli, F., Ranzuglia, G.: Meshlab: an open-source mesh processing tool. In: *Eurographics* (2008)
- [13] Clark, J.J.: Active photometric stereo. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (1992)
- [14] Collins, T., Bartoli, A.: 3D Reconstruction in Laparoscopy with Close-Range Photometric Stereo. In: *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. Nice, France (2012)
- [15] Deguchi, K., Okatani, T.: Shape reconstruction from an endoscope image shape-from-shading technique for a point light source at the projection center. In: *Workshop on MMBIA* (1996)
- [16] Frankot, R.T., Chellappa, R.: A method for enforcing integrability in shape from shading algorithms. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)* **10**(4), 439–451 (1988)
- [17] G, C., K., H., B., S., Y., M., K., W.K.Y.: Deep photometric stereo for non-Lambertian surfaces. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)* (2020)
- [18] Hinton, G.E.: Deep belief networks. *Scholarpedia* (2009)
- [19] Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2017)
- [20] Hui, Z., Sankaranarayanan, A.C.: Shape and spatially-varying reflectance estimation from virtual exemplars. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2016)

- [21] Ikehata, S.: Cnn-ps: Cnn-based photometric stereo for general non-convex surfaces. In: European Conference on Computer Vision (ECCV) (2018)
- [22] Ikehata, S., Aizawa, K.: Photometric stereo using constrained bivariate regression for general isotropic surfaces. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2014)
- [23] Ikehata, S., Wipf, D., Matsushita, Y., Aizawa, K.: Robust photometric stereo using sparse regression. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2012)
- [24] Iwahori, Y., Sugie, H., Ishii, N.: Reconstructing shape from shading images under point light source illumination. In: ICPR (1990)
- [25] Ju, Y., Qi, L., Zhou, H., Dong, J., Lu, L.: Demultiplexing colored images for multispectral photometric stereo via deep neural networks. IEEE Access (2018)
- [26] Konstantinou, C., Biscontin, G., Logothetis, F.: Tensile strength of artificially cemented sandstone generated via microbially induced carbonate precipitation. *Materials* **14**(16), 4735 (2021)
- [27] Lee, S., Brady, M.: Integrating stereo and photometric stereo to monitor the development of glaucoma. *Image and Vision Computing* (1991)
- [28] Li, Z., Xu, Z., Ramamoorthi, R., Sunkavalli, K., Chandraker, M.: Learning to reconstruct shape and spatially-varying reflectance from a single image. *ACM Transactions on Graphics (TOG)* **37**(6), 1–11 (2018)
- [29] Liu, C., Narasimhan, S.G., Dubrawski, A.W.: Near-light photometric stereo using circularly placed point light sources. In: ICCV (2018)
- [30] Logothetis, F., Budvytis, I., Mecca, R., Cipolla, R.: A cnn based approach for the near-field photometric stereo problem. In: British Machine Vision Conference (BMVC) (2020)
- [31] Logothetis, F., Budvytis, I., Mecca, R., Cipolla, R.: PX-NET: Simple, Efficient Pixel-Wise Training of Photometric Stereo Networks. In: ICCV (2021)
- [32] Logothetis, F., Mecca, R., Cipolla, R.: Semi-calibrated near field photometric stereo. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
- [33] Logothetis, F., Mecca, R., Cipolla, R.: A differential volumetric approach to multi-view photometric stereo. In: ICCV (2019)
- [34] Logothetis, F., Mecca, R., Quéau, Y., Cipolla, R.: Near-field photometric stereo in ambient light. In: British Machine Vision Conference (BMVC) (2016)
- [35] Matusik, W., Pfister, H., Brand, M., McMillan, L.: A data-driven reflectance model. *ACM Transactions on Graphics* (2003)
- [36] Mecca, R., Falcone, M.: Uniqueness and approximation of a photometric shape-from-shading model. *SIAM J. Imag. Sci.* **6**(1), 616–659 (2013)
- [37] Mecca, R., Logothetis, F., Budvytis, I., Cipolla, R.: Lucas: A dataset for near-field point light source photometric stereo. In: British Machine Vision Conference (BMVC) (2021)
- [38] Mecca, R., Quéau, Y., Logothetis, F., Cipolla, R.: A single lobe photometric stereo approach for heterogeneous material. *SIAM Journal on Imaging Sciences* (2016)
- [39] Mecca, R., Rodolà, E., Cremers, D.: Realistic photometric stereo using partial differential irradiance equation ratios. *Computers & Graphics* **51**, 8–16 (2015)
- [40] Mecca, R., Wetzler, A., Bruckstein, A., Kimmel, R.: Near Field Photometric Stereo with Point Light Sources. *SIAM Journal on Imaging Sciences* (2014)
- [41] Onn, R., Bruckstein, A.: Integrability disambiguates surface recovery in two-image photometric stereo. *IJCV* (1990)

- [42] Parot, V., Lim, D., González, G., Traverso, G., Nishioka, N., Vakoc B.J. and Durr, N.: Photometric stereo endoscopy. *J Biomed Opt.* (2013)
- [43] Prados, E., Faugeras, O.: Perspective shape from shading and viscosity solutions. In: *ICCV* (2003)
- [44] Quéau, Y., Durix, B., Wu, T., Cremers, D., Lauze, F., Durou, J.D.: Led-based photometric stereo: Modeling, calibration and numerical solution. *JMIV* (2018)
- [45] Quéau, Y., Durou, J.D.: Edge-preserving integration of a normal field: Weighted least squares, TV and L1 approaches. In: *SSVM* (2015)
- [46] Quéau, Y., Mecca, R., Durou, J.D.: Unbiased photometric stereo for colored surfaces: A variational approach. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016)
- [47] Quéau, Y., Modrzejewski, R., Gurdjos, P., Durou, J.D.: A Full Photometric and Geometric Model for Webcam + Matte Screen Devices. *Signal Processing: Image Communications* **40**, 65–81 (2016)
- [48] Saiz, F.A., Barandiaran, I., Arbelaz, A., Graña, M.: Photometric stereo-based defect detection system for steel components manufacturing using a deep segmentation network. *Sensors* (2022)
- [49] Santo, H., Samejima, M., Sugano, Y., Shi, B., Matsushita, Y.: Deep photometric stereo network. In: *ICCV Workshops* (2017)
- [50] Santo, H., Waechter, M., Matsushita, Y.: Deep near-light photometric stereo for spatially varying reflectances. In: *European Conference on Computer Vision (ECCV)* (2020)
- [51] Shi, B., Mo, Z., Wu, Z., Duan, D., Yeung, S.K., Tan, P.: A benchmark dataset and evaluation for non-lambertian and uncalibrated photometric stereo. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2018)
- [52] Shi, B., Wu, Z., Mo, Z., Duan, D., Yeung, S.K., Tan, P.: A benchmark dataset and evaluation for non-lambertian and uncalibrated photometric stereo. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016)
- [53] Simchony, T., Chellappa, R., Shao, M.: Direct analytical methods for solving poisson equations in computer vision problems. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)* **12**(5), 435–446 (1990)
- [54] Smith, W., F., F.: Height from photometric ratio with model-based light source selection. *CVIU* (2016)
- [55] Tang, Y., Salakhutdinov, R., Hinton, G.E.: Deep lambertian networks. In: *ICML* (2012)
- [56] Taniai, T., Maehara, T.: Neural inverse rendering for general reflectance photometric stereo. In: *ICML* (2018)
- [57] Tankus, A., Kiryati, N.: Photometric stereo under perspective projection. In: *ICCV* (2005)
- [58] Vogiatzis, G., Hernández, C.: Practical 3d reconstruction based on photometric stereo. In: *Computer Vision: Detection, Recognition and Reconstruction*. Springer (2010)
- [59] Wang, X., Jian, Z., Ren, M.: Non-lambertian photometric stereo network based on inverse reflectance model with collocated light. *IEEE Transactions on Image Processing: a Publication of the IEEE Signal Processing Society* (2020)
- [60] Wetzler, A., Mecca, R., Bruckstein, A.M., Kimmel, R.: Close-range photometric stereo with point light sources. In: *3DV* (2014)
- [61] Wolff, L.B., Angelopoulou, E.: 3-D stereo using photometric ratios. In: *European Conference on Computer Vision (ECCV)*, pp. 247–258 (1994)
- [62] Woodham, R.J.: Photometric method for determining surface orientation from multiple images. *Optical Engineering* (1980)

- [63] Wu, C., Narasimhan, S.G., Jaramaz, B.: A Multi-Image Shape-from-Shading Framework for Near-Lighting Perspective Endoscopes. *International Journal of Computer Vision (IJCV)* (2010)
- [64] Xiong, Y., Chakrabarti, A., Basri, R., Gortler, S.J., Jacobs, D.W., Zickler, T.E.: From shading to local shape. *IEEE Trans. Pattern Anal. Mach. Intell.* **37**(1), 67–79 (2015)
- [65] Yu, Y., Smith, W.A.P.: Pvnnet: A neural network library for photometric vision. In: *ICCV Workshop* (2017)
- [66] Yuille, A.L., Snow, D., Epstein, R., Belhumeur, P.N.: Determining generative models of objects under varying illumination: Shape and albedo from multiple images using SVD and integrability. *IJCV* (1999)
- [67] Zheng, Q., Jia, Y., Shi, B., Jiang, X., Duan, L.Y., Kot, A.C.: Spline-net: Sparse photometric stereo through lighting interpolation and normal estimation networks. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2019)
- [68] Zhu, D., Smith, W.A.P.: Least squares surface reconstruction on arbitrary domains. In: *European Conference on Computer Vision (ECCV)* (2020)