

DiLiGenT-II: Photometric Stereo for Planar Surfaces with Rich Details – Benchmark Dataset and Beyond

Feishi Wang^{1,2,3,†} Jiejie Ren^{4,†} Heng Guo^{5,6,†} Mingjun Ren^{4,‡} Boxin Shi^{1,2,3,‡}

¹ National Key Laboratory for Multimedia Information Processing, School of Computer Science, Peking University

² National Engineering Research Center of Visual Technology, School of Computer Science, Peking University

³ AI Innovation Center, School of Computer Science, Peking University

⁴ School of Mechanical Engineering, Shanghai Jiao Tong University

⁵ School of Artificial Intelligence, Beijing University of Posts and Telecommunications ⁶ Osaka University

{wangfeishi, shiboxin}@pku.edu.cn, {jiejieren, renmj}@sjtu.edu.cn, guoheng.bupt@gmail.com

Abstract

Photometric stereo aims to recover detailed surface shapes from images captured under varying illuminations. However, existing real-world datasets primarily focus on evaluating photometric stereo for general non-Lambertian reflectances and feature bulgy shapes that have a certain height. As shape detail recovery is the key strength of photometric stereo over other 3D reconstruction techniques, and the near-planar surfaces widely exist in cultural relics and manufacturing workpieces, we present a new real-world dataset DiLiGenT-II containing 30 near-planar scenes with rich surface details. This dataset enables us to evaluate recent photometric stereo methods specifically for their ability to estimate shape details under diverse materials and to identify open problems such as near-planar surface normal estimation from uncalibrated photometric stereo and surface detail recovery for translucent materials. To inspire future research, this dataset will open soruced at <https://photometricstereo.github.io/diligentpi.html>.

1. Introduction

Photometric stereo [47, 44] aims at single view three-dimensional (3D) reconstruction from image observations captured under varying lights. Compared to structured light-based 3D reconstruction techniques that are widely applied in commercial scanners, the key strength of photometric stereo is the detailed surface shape recovery, which is of great interest for additive manufacturing and rendering.

To evaluate the effectiveness of photometric stereo methods, a batch of real-world benchmark datasets have been

built such as DiLiGenT [43, 41] (and its multi-view extension DiLiGenT-MV [29]), DiLiGenT10² [39] for distant lights, and LUCES [34] for near lights. Existing datasets focus on evaluating the effectiveness of photometric stereo techniques on non-Lambertian surfaces, revealing the performance of existing methods on real-world scenes. Since the majority of target objects in these datasets are smooth surfaces (some of them with a portion of detailed structures on the surface), it is hard to evaluate the accuracy of surface detail recovery, which is unique to photometric stereo over other 3D reconstruction techniques. On the other hand, near-planar surfaces usually accompanied with rich details are commonly observed in our daily life, such as reliefs, badges, and coins. However, existing photometric stereo datasets mainly choose statue-like objects or other bulgy shapes with a certain height as targets, lacking sophisticated evaluation on near-planar surfaces.

In this paper, we build a new dataset named **DiLiGenT-II**¹ for photometric stereo focusing on the recovery of near-planar shape and surface details. As shown in Fig. 1, we collect 30 representative real-world planar objects with rapidly varied geometric details. The dataset can be categorized into 4 groups containing **metallic**, **specular**, **rough**, and **translucent** surface reflectance, respectively. The metallic group contains 10 metallic coins; the specular group includes 10 enamel badges; the translucent group has 5 3D-printed objects made by photo-polymer resin, and the rough group contains 5 surfaces sharing the same geometry of the translucent group but sprayed with a matte paint [35]. In addition, our DiLiGenT-II takes an optical profilometer to capture the ultra-precise surface 3D structure in nanometer

¹ ‘DiLiGenT’ [41] as the abbreviation of Directional Lighting, General reflectance, with the ‘ground Truth’ shapes for photometric stereo benchmarking. As we take the similar assumptions, we refer DiLiGenT as prefix and use II to indicate planar objects.

[†] Equally contributed authors [‡] Corresponding authors.



Figure 1: Overview of DiLiGenT-II. We collect 4 groups of near-planar objects with rich surface details and diverse reflectance types (metallic, specular, rough, and translucent). The corresponding ‘ground-truth’ surface normals shown in the even rows are measured via a precise profilometer in the accuracy of nanometers.

accuracy, providing the ‘ground truth’ surface normal with well-preserved tiny surface details.

We apply DiLiGenT-II to evaluate up-to-date photometric stereo methods under the settings of calibrated and uncalibrated distant light, and benchmark the reconstruction performance on detailed structures and near-planar surfaces. The evaluation results reveal the difference in surface detail recovery of learning-based photometric stereo methods working with per-pixel and all-pixel manners [54]; the challenging problems of uncalibrated photometric stereo for near-planar surfaces; and the influence of translucent and rough reflectance on surface detail recovery. The analysis on the DiLiGenT-II presents new challenges and open problems for photometric stereo.

To summarize, this paper contributes to photometric stereo benchmark and inspires future research by proposing:

- the first real-world dataset, DiLiGenT-II, that evaluates

- near-planar surfaces with rich geometric details;
- up-to-date benchmark evaluation of photometric stereo recovering important features for handling surface details; while
- revealing inherent obstacles of planar detailed objects to photometric stereo with open problems.








2. Related Work

This paper focuses on the benchmark dataset for evaluating photometric stereo on near-planar surfaces with rich details. In the following sections, we will discuss the related datasets and representative works in photometric stereo.

2.1. Photometric Stereo Datasets

Synthetic datasets adopt physical-based rendering engines such as Mitsuba [22] and Blender [11] to create image observations and the corresponding surface normal maps of

Table 1: Summary of real-word photometric stereo datasets. Material: controlled (C: fabricated or carefully selected with controlled categories) or uncontrolled (UC: randomly picked up from daily objects); Ground Truth (GT measure): from CAD/Scanned models with registration (+Reg) or from photometric stereo (PS). Number (#) of shapes, lights, and sets (one set means a sequence of photometric stereo images under varying lighting conditions used for computation).

Dataset							
	DiLiGenT-II	DiLiGenT10 ² [39]	DiLiGenT [41]	LUCES [34]	Harvard [51]	ETHz [27]	Gourd&Apple [4]
GT measure	Scan+Reg	CAD+Reg	Scan+Reg	Scan+Reg	PS	Scan+Reg	PS
Material	C	C	UC	UC	UC	UC	UC
# Shapes	30	10	10	14	7	3	2
# Lights	100	100	96	52	20	260	102/112
# Sets	30	100	10	14	7	3	2

diverse synthetic scenes. DPSN [40] introduces the first synthetic photometric stereo dataset BlobbyPS, comprising 10 smooth Blobby [23] shapes with reflectance assigned by measured MERL BRDFs [33] rendered under 96 light directions. To extend the smooth surface of Blobby to more complex shapes, PS-FCN [8] takes 59, 292 scanned sculptures to render the SculpturePS dataset. CNN-PS [18] proposes the CyclePS dataset to extend material distribution from uniform to spatially-varying, where each sub-region or even each pixel of the surface is assigned with distinct reflectances modeled by Disney Principle BSDF [5]. While synthetic datasets expand photometric stereo data availability, existing synthetic datasets are mostly rendered with statue-like objects such as the shapes in Blobby [23] and Sculpture dataset [6], leading to a shape domain gap to flatten surfaces. In this paper, we render a synthetic dataset containing 127 near-planar surfaces with rich details to enhance learning-based photometric stereo on near-planar surface recovery.

Real-world datasets complement the gap between computer graphics rendering and real-world imaging process. The Gourd&Apple dataset [2] releases image observations of 2 objects with spatially-varying isotropic BRDFs. Harvard dataset [50] contains 7 surfaces with uniform diffuse reflectance. However, the ground truth surface normal of the above two datasets is not provided. DiLiGenT [41] records 10 objects with diverse shapes and general non-Lambertian materials. Starting from DiLiGenT [41], benchmark evaluation of photometric stereo becomes available based on the ‘ground truth’ surface normal from scanned meshes. The following-up datasets further extends DiLiGenT [41] from the perspective of multi-views (DiLiGenT-MV [30]), controlled materials and shapes (DiLiGenT 10² [39]), near-field illumination (LUCES [34]), environment illumination [1, 16, 17], and global illumination effects [27].

As summarized in Tab. 1, objects contained in existing real-world photometric stereo datasets are mostly bulgy and smooth, which do not include flattened or near-planar surfaces though they are commonly seen in our daily life. More importantly, the smooth target shapes are not suitable to evaluate the uniqueness of photometric stereo over other 3D reconstruction techniques, that is, the high-fidelity surface detail recovery, especially for delicate structures. To address these two problems, we newly build a real-world photometric stereo dataset containing near-planar surfaces with rich geometric details.

2.2. Photometric Stereo Methods

We briefly review existing photometric stereo methods based on distant calibrated and uncalibrated light settings. Please refer to the photometric stereo survey of non-learning based methods [41] and learning-based methods [54, 25] for more comprehensive analysis.

Calibrated photometric stereo works with calibrated distant lights. Existing non-learning based methods either treat specular highlights and shadows as sparse outliers [48, 38] or propose parametric or data-based reflectance model to explicitly handle the non-Lambertian reflection [10, 14, 42, 15]. Beginning from DPSN [40], the image observations are directly mapped to the corresponding surface normal via deep neural networks, where the non-Lambertian reflectance is implicitly learned from the synthetic training data. The network structure in existing learning-based methods can be divided into the all-pixel branch (representative work: PS-FCN [8]) and the per-pixel branch (representative work: CNN-PS [18]). Based on these two typical network structures, following-up works further promote photometric stereo by addressing the global illumination effects (e.g. PX-Net [32]), reducing the number of inputs (e.g. SPLINE-Net [53], LMPS [28]), combining the

merit of per-pixel and all-pixel methods (e.g. GPS-Net [52], PS-Transformer [19]). However, few methods focus on the recovery of surface details despite its importance in photometric stereo. Ju *et al.* [24] firstly addressed the high-frequency surface details in photometric stereo based on attention-weighted loss. Their following-up work [26] further enhanced the detail recovery accuracy via a double-gate normalization and a parallel high-resolution structure.

Uncalibrated photometric stereo works under unknown distant light directions, so that photometric stereo becomes more challenging even with Lambertian reflectance assumption. Non-learning based methods adopt the local diffuse reflectance (LDR) maxima [37], perspective camera [36], specularities [13], albedo entropy [3] to resolve the surface normal estimation ambiguity in uncalibrated Lambertian photometric stereo. Beginning from SDPS-Net [7], learning-based uncalibrated photometric stereo methods capable of handling non-Lambertian reflectance are proposed. SDPS-Net [7] first estimates the light directions and intensities from input image observations, and then feeds them into the normal estimation network. Kaya *et al.* [27] proposed an uncalibrated deep neural network with an inverse rendering module, where the inter-reflections are explicitly modeled in the image formation process. Following the analysis of Chen *et al.* [9], the light calibration in deep uncalibrated photometric stereo is related to the existence of attached shadows and specular highlights in the image observations. For near-planar surfaces where surfaces are flattened and attached shadows are rarely observed, whether the existing methods can be applied is not evaluated. The above uncalibrated photometric stereo methods assume distant light. Existing methods further extend the illumination setting to general natural light. Lichy *et al.* [31] proposes a weakly calibrated method working in the indoor environment, requiring at most 6 images from approximately known lighting directions. Ikehata proposes photometric stereo methods [20, 19] to handle general unknown illumination in real-world scenes.

3. DiLiGenT-II Dataset

This section introduces our DiLiGenT-II dataset, which contains 30 near-planar scenes covered by varying materials and geometric details. Each real-world scene contains RGB images under 100 varying light directions and a precisely measured ‘ground-truth’ normal map. The resolution for each scene is 960×960 .

3.1. Objects groups

As shown in Fig. 1, we collect 4 groups of near-planar objects named by their reflectance properties: **metallic**, **specular**, **rough**, and **translucent**. In each group, we select target objects with rich geometric details with size around

15 ~ 25 mm.

Metallic group. We choose 10 different coins as target objects commonly observed in our daily life. These coins are casted from different metallic materials, such as nickel, brass, aluminum alloy, and bi-metal. The reflectance distributions are either uniform (e.g. CRAB) or spatially-varying (e.g. RHINO). The surface normal of the coins contains rich and delicate geometric details and even traces of daily use. Their corresponding surface PV (peak to valley) depth value is usually less than 1 mm.

Specular group. We choose 10 sets of badges as the target objects. These badges are made of polished metal, plastic, and enamel, showing strong and sparse specular spikes (e.g. BEAR), or broad and soft specular lobes (e.g. TREE) on their captured images. All objects in this category contain spatially-varying reflectances and have greater depth variation (PV is around 1.5 mm), making this group different from the metallic group.

Translucent group. We collect 5 sets of 3D-printed relief surfaces to build the translucent object group. The materials for 3D printing for the LOTUS-T and the BAGUA-T are photo-polymer resin² with different colors, respectively, which is known to be slightly translucent, bringing unavoidable subsurface scattering in the surface reflectance. The surface shapes in this group have even stronger surface undulating (PV is around 4 mm), leading to stronger shadows and inter-reflections being observed in the captured images.

Rough group. We fabricate another 5 sets of 3D-printed relief surfaces sharing exactly the same shapes with the translucent group but covered by gray matte spray³, whose reflectance is approximately close to the Oren-Nayar diffuse model [35] and the Torrance-Sparrow specular model [46], as analyzed in [45].

3.2. Capture System and Light Calibration

As shown in Fig. 2, we design and build a photometric stereo image capture setup with the function of automatic illumination and capture at varying light directions. The capture system is placed into a darkroom cage covered with black felt to shield the ambient lights.

On the illumination side, we build a dual-axis rotation platform to control the light direction, which is based on robot arms capable of omni-direction illuminations on the upper hemisphere of the target surface. We attach a single LED light source on our robot arm to ensure the consistency of emitted light intensity among different images. We further adopt a co-concentric rotation design to achieve

²Kexcelled Ultradetail: <https://www.kexcelled3d.com/products/ultradetail/>. Retrieved August 16, 2023.

³FA-5: <http://cysygroup.com/en/product.asp?category=NDT&page=5>. Retrieved August 16, 2023.

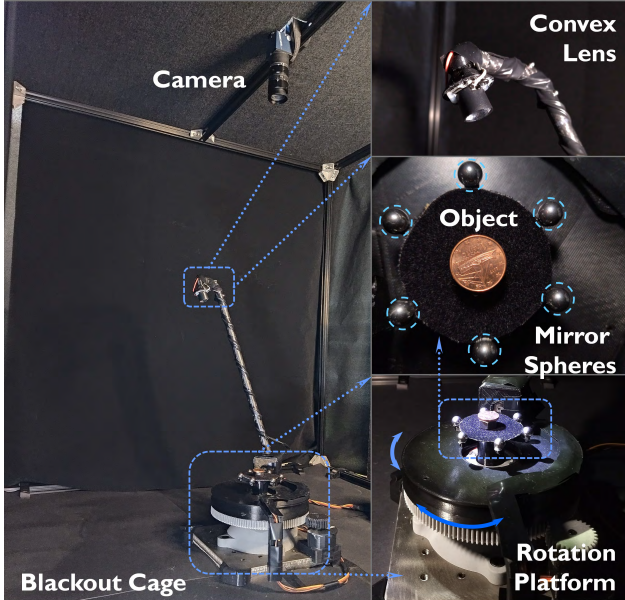


Figure 2: Capture system overview. Our automatic photometric stereo image capture equipment is placed in the blackout cage, containing two rotation axes to provide illumination from an arbitrary direction from the upper hemisphere. A lens is placed in front of the LED to improve the directionality and uniformity. Six mirror balls around the object are used to calibrate light directions ‘on the fly’.



Figure 3: Process of obtaining a ‘GT’ normal map.

roughly the same distance between light and the target surface so that the received light intensities at different surface positions are roughly the same. To improve the directivity of the LED point light source, A convex lens is placed in front of the LED source.

On the camera side, we use a Daheng MER-503-36U3C camera⁴ with a 50 mm lens to record 12-bit raw images (with linear radiometric response). Under each directional light, we capture 10 images with the exposure time settings from 1 ms to 10 ms, from which a high dynamic range (HDR) image that records both dark shadows and

⁴Camera website: <https://www.get-cameras.com/USB-Camera-Sony-IMX264-MER-503-36U3C>. Retrieved August 16, 2023.

bright specular highlights can be reconstructed based on the existing HDR algorithm [12].

During the capture, the object is placed in the center of the rotation platform. The system controls the robot arm to move toward a pre-defined light direction and triggers the camera to shoot at various exposure times to obtain an HDR image measurement. Then, the robot arm is rotated to the next light direction. Totally, 100 HDR images under varying lights are recorded by repeating the above process.

To calibrate the light directions, we place 6 mirror balls with a radius of 8 mm around the object stage, and record the specular highlights at each mirror ball during data capture. The light directions can then be calculated based on the specular positions following the same calibration method described in existing methods [41, 39]. Please check our supplementary material for more details.

3.3. Obtaining the ‘GT’ Normal Map

Similar to previous real-world datasets: DiLiGenT [41] and LUCES [34], we measure the ‘ground truth’ surface normal from scanned mesh in DiLiGenT-II. As shown in Fig. 3, we apply a commercial 3D scanner Bruker Alicona Infinity Focus (up to **10 nm** accuracy)⁵ to measure a precise point cloud. The scanner is based on Focus-Variation measurement and can probe vertical surfaces precisely, which is necessary to scan near-planar objects completely.

Given a scanned mesh of the near-planar surface, we manually adjust the camera pose to align the mesh with one captured image by matching key points and geometric features. Then, taking the calibrated intrinsic camera parameters and the extrinsic pose, we render the mesh to the corresponding surface normal map by Blender [11] with the same resolution to the captured images. We try our best to check key points accuracy at the sub-pixel level, but inevitable errors in the manual alignment process might still exist, so we add a quotation on the ‘ground truth’ like [41, 39].

4. Benchmark Analysis

This section showcases the benchmark results for photometric stereo techniques using the DiLiGenT-II dataset.

4.1. Baseline Methods & Evaluation Metric

Based on the survey [41] adopt non-Learning photometric stereo, we choose the baseline method (least-square based Lambertian photometric stereo [47], LSPS), baseline method with position thresholding strategy [41] (TH28 and TH46 reject pixels whose intensities under varying lights are outside the range of [20%, 80%] and [40%, 60%]), respectively), WG10 [49] (a robust photometric stereo method based on outlier rejection), ST14 [42] (showing

⁵Alicona website: <https://www.alicona.com/en/products/infinitefocus>. Retrieved August 16, 2023.

Table 2: Benchmark results on our real-world dataset DiLiGenT-II. Mean angular error in degree of each object on various methods are presented, and the average angular error for each material group and for all objects is shown in the bottom rows. We denote ‘NA-PSN’ as the abbreviation of NormAttention-PSN [26].

		Calibrated Methods										Uncalibrated Methods				
	Plane	LSPS [47]	TH28 [41]	TH46 [41]	WG10 [49]	ST14 [42]	PX- Net [32]	CNN- PS [18]	PS- FCN [8]	GPS- Net [52]	NA- PSN [26]	PF14 [37]	UPS- FCN [8]	SDPS- Net [7]	LW21 [31]	SDM- UniPS [20]
Metallic	FLOWER	5.8	6.8	7.5	6	7.2	7.9	5.5	4.7	4.6	4.6	35.3	12.7	12.8	23.9	15.2
	BIRD	9	8.1	9.1	8	11.6	10.8	8.8	6.8	7.2	6.8	34.9	15.4	17.6	28.5	26.5
	RHINO	6.8	7.5	7.8	8.9	9.3	8.9	10.6	4.9	5.3	5.6	30.2	17.7	24.9	26.8	17
	LIONS	6.6	6.7	7.7	7.4	9.1	8.1	6.8	4.7	4.5	4.6	40.9	11.8	19.6	27.1	9.2
	QUEEN	6.2	7.5	8.1	8.9	8.6	8.3	14.7	5.4	4.7	4.7	46.4	12.6	16.5	27.2	10.6
	CRAB	5.9	7.1	7.6	8.6	8.5	8.4	7.2	4.5	5.3	4.9	31.7	18.1	20.5	26.8	25.4
	SHIP	5.9	7.1	8.9	6.4	9.3	9.1	12.3	4.9	6.1	5.1	89.7	15.6	19	26.4	22
	PARA	6.7	6	6.6	5.1	8.4	8	4.9	3.9	4.7	4	42.7	17	19.8	24.6	23.2
	SAIL	6.9	8.7	10.2	8.8	8.1	8.2	13.4	5.2	5.1	5.5	40.6	13.5	16.7	27	10.5
	FISH	5.2	6.7	7.9	6.5	8.2	8.1	7.4	4.2	4.6	4.6	39.9	15.8	23.6	25.6	24.5
Specular	TREE	20.1	9.4	11	9.7	14.9	10.6	11.3	7.8	10.2	9.3	32	40.6	34.1	35.8	47.2
	OCEAN	14.5	6.2	6	5.9	8.4	7.6	5.9	4.6	5.8	5.4	77.9	26.7	31.4	28.1	34.6
	LUNG	20.3	8.2	8.8	8.1	11.8	9.4	7.6	5.7	9.7	7.5	41.5	36.5	40.2	34.4	46.6
	BEAR	9.8	9.1	8.5	8.9	8.6	9.7	7.9	7.4	7.4	6.9	84.8	27.6	30.7	27	23.8
	TV	19.6	11.7	13.3	13.2	17.1	13.5	12.1	11.3	10.6	10.6	73	22.8	41.1	35.4	34.4
	SUN	11.9	7.9	9.5	11.7	10.5	8.2	6.7	5.8	6.7	8	27.9	22.8	31.5	24	26.2
	TAICHI	17.3	9.2	11.2	11	17.9	10.7	8.7	8.3	8	8.3	58.9	29.8	26.9	28.9	36.6
	WAVE	15.7	6.9	7.7	7.3	9.2	8.4	7.1	5.3	6.8	6.3	76	30.8	39.1	25.9	34.9
	ASTRO	17.6	7.9	8.7	7.9	11.4	9.2	7.1	6	7.2	7.7	35.2	25.9	37.7	36.5	37.8
	WHALE	16.3	8.8	8.9	8.5	13.4	10.1	9.5	11.6	12.2	10.4	63.2	37.6	29.8	32	33.8
Translucent	BAGUA-T	20.6	16.7	16.5	16.1	17.8	17.4	15.8	16.4	16.8	16.1	20.3	19.9	28.9	27.8	17.1
	LOTUS-T	15.9	14	14.1	13.8	14.3	14.5	13.7	13.5	13.6	13.7	19.8	19.5	26.5	21.3	13.6
	LION-T	30.1	21	20.5	19.9	21.6	21.3	18.4	20.3	21.2	23.4	22.7	25.8	23.6	30.5	16.2
	PANDA-T	18.6	16.4	16.4	16.3	16.8	16.8	15.9	16.6	17.2	17	19.2	18.5	23.7	22.8	17.6
	CLOUD-T	18.7	17	17	16.8	18.1	17.7	16.1	17.2	17.8	17.6	20.6	20.7	27.5	24.2	19.2
Rough	BAGUA-R	20.4	13.2	12.2	11.6	16.5	13.6	11.9	12.2	13	16.4	19.5	16.8	22.5	25	14.6
	LOTUS-R	15.7	12	11.7	11.3	12.3	12.4	11	10.9	11.8	13.4	18.3	14.7	21.7	23.8	11.8
	LION-R	30.2	19.4	17.8	17	19.4	19.1	17.9	15.8	18.4	23	19.5	25.2	20.8	29.7	15.9
	PANDA-R	18.7	14.6	14.1	14	16	15	14.1	14.2	14.8	16.3	17.4	16.8	21.8	21.6	17.1
	CLOUD-R	18.4	14	13.6	13.5	17.4	14.7	13.6	14.6	14.3	15.9	18.4	16.8	27.4	21.7	17.1
M.Avg.		6.5	7.2	8.1	7.5	8.8	8.6	9.2	4.9	5.2	5.1	43.2	15	19.1	26.4	18.4
S.Avg.		16.3	8.5	9.4	9.2	12.3	9.7	8.4	7.4	8.5	8	57	30.1	34.3	30.8	35.6
T.Avg.		20.8	17	16.9	16.6	17.7	17.5	16	16.8	17.3	17.6	20.5	20.9	26.1	25.3	16.8
R.Avg.		20.7	14.6	13.9	13.5	16.3	15	13.7	13.5	14.4	17	18.6	18.1	22.8	24.4	15.3
Avg.		14.5	10.5	11	10.6	12.7	11.5	10.8	9.2	9.8	10.1	39.9	21.5	25.9	27.3	23.3

best performance at DiLiGenT reported in [41]). Based on surveys about the learning-based photometric stereo [54, 25], we choose representative networks handling photometric stereo in an all-pixel manner (NormAttention-PSN [26], PS-FCN [8]), per-pixel manner (CNN-PS [18], PX-Net [32]), and the method GPS-Net [52] that considers both the per-pixel and the all-pixel structures. Besides the above calibrated photometric stereo methods, we also eval-

uate existing uncalibrated photometric stereo approaches, including a non-learning based method PF14 [37] (showing the best performance under uncalibrated setting as reported in [41]), three representative learning-based methods SDPS-Net [7], UPS-FCN [8] and UPS-GCNet [9] published in recent years. We also include photometric stereo methods such as LW21 [31] and SDM-UniPS [20] that are based on more general natural illumination setup. During the evalu-

ation, we adopted the code and pre-trained model released by the authors to process the collected data in DiLiGenT-II.

Similar to DiLiGenT [41] and DiLiGenT10² [39], the mean angular error (MAngE) between the estimated and the ground-truth surface normal is used as the metric for measuring the performance of photometric stereo methods quantitatively.

To clarify any concerns regarding the benchmark results, it’s important to note that the input format varies between optimization-based methods (LSPS [47], TH28, TH46, WG10 [49], ST14 [42], and PF14 [37]) and other learning-based methods. Optimization-based methods employ a float format, providing higher precision, while other learning-based methods utilize quantized 16-bit input. This differentiation arises because the optimization-based methods are not sensitive to changes in the pixel value range resulting from quantization. For the learning-based methods, adhering to the standard procedure and employing quantized images as input ensures the reliability of the evaluation results.

4.2. Analysis to Different Baselines

As shown in Table 2, we evaluate the surface normal estimation accuracy for all methods on DiLiGenT-II. In the following, we first analyze the surface detail recovery from calibrated photometric stereo, followed by the analysis of uncalibrated photometric stereo on near-planar surfaces.

Calibrated photometric stereo. From Table 2, we observe that the LSPS [47] shows the best performance over other non-learning-based photometric stereo methods, while the angular error difference between LSPS [47], TH28, TH46, and ST14 [42] are marginal. For surfaces with shadows such as LION_R, or containing dominant Lambertian reflectances and sparse specular highlights such as FLOWER, TH46 is more effective than the LSPS [47]. However, TH46 could be unstable as only 20% of the image observations under varying lights are used for computing surface normals, especially for surfaces with dark reflectances (e.g.SUN and Tv). These observations are consistent with the previous evaluation on DiLiGenT10² [39].

Among learning-based calibrated photometric stereo, CNN-PS [18] and NormAttention-PSN [26] achieve smaller mean angular errors, showing their advantages on detailed surface recovery. As shown in Fig. 4, we visualize the error distributions of PS-FCN [8], CNN-PS [21], and NormAttention-PSN [26] on four representative surfaces belong to the four reflectance groups. PS-FCN [7] outputs blurry normal estimation results as highlighted in Fig. 4, possibly due to the spatial smoothness brought by the fully convolutional networks. As CNN-PS [18] solves photometric stereo in a per-pixel manner, the surface details are not contaminated by the neighboring pixels. NormAttention-PSN [26] is built upon PS-FCN [8] but can handle surface

detail recovery by an attention-weighted loss. We summarize the benchmark results of calibrated photometric stereo using DiLiGenT-II as the following observation:

Observation 1 *Learning-based photometric stereo methods in the per-pixel branch are generally more effective than the all-pixel branch for handling surface detail estimation. A detail-weighted loss can help methods in the all-pixel branch for better recovering tiny structures.*

Uncalibrated photometric stereo. From Table 2, uncalibrated photometric stereo methods generally show significant errors on near-planar surfaces in DiLiGenT-II compared to the case of calibrated photometric stereo. For the non-learning-based method PF14 [37], the near-planar surface could be an ill-posed shape when conducting SVD for obtaining pseudo-surface normals and lights. Also, for learning-based uncalibrated photometric stereo, attached shadows and shading variations are essential in recovering the unknown light directions, as stated in [9]. However, the shadows are much less observable for near-planar surfaces as the tiny detail can barely cast a block of shadows from varying light directions compared with the case of the bulgy objects, making the light estimation more challenging.

Fine-tuning on the synthetic near-planar dataset. To check whether the error of learning-based uncalibrated photometric stereo can be further reduced by learning the data prior, we create a near-planar synthetic dataset PS_RELIEF for finetuning the uncalibrated photometric stereo. As shown in Fig. 5, our synthetic dataset contains 127 near-planar surface normals extracted from CAD meshes. We adopt Disney’s principled BSDF [5] as the reflectance model and randomly generate BRDFs by adjusting the parameter to control the diffuse, specular, and metallic reflectance components in a similar manner to the existing synthetic dataset CyclePS [18]. In total, PS_RELIEF contains 3429 scenes. For each scene, we render the object by Blender [11] under 100 distant light directions with the same uniform lighting distribution of DiLiGenT10² [39].

Table 3: Ablation study on photometric stereo methods trained with or without fine-tuning (FT.) on the near-planar synthetic dataset PS_RELIEF.

Light	Method	DiLiGenT [41]		DiLiGenT-II	
		w/ FT.	w/o FT.	w/ FT.	w/o FT.
Uncalibrated	SDPS-Net [7]	17.80	9.51	16.03	27.76
	UPS-Net [8]	25.05	15.37	14.99	21.54
Calibrated	PS-FCN [8]	11.75	9.12	8.80	9.84

As shown in Table 3, we fine-tune the UPS-FCN [8] and SDPS-Net [7] with PS_RELIEF and test it on real-world dataset DiLiGenT [41] and DiLiGenT-II, which contain bulgy shapes and near-planar shapes, respectively. The

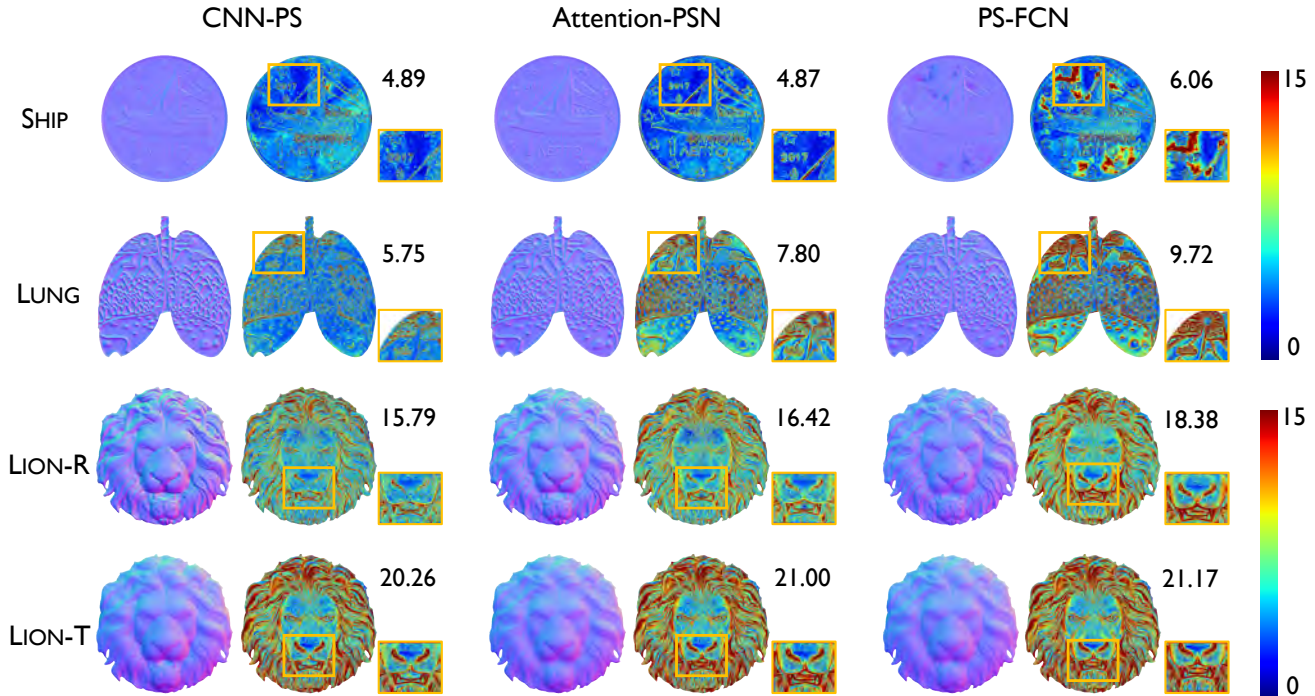


Figure 4: Visualization of surface detail recovery from different photometric stereo methods, where the odd and even columns plot the estimated surface normals and the corresponding angular error maps, respectively. Per-pixel based photometric stereo method CNN-PS [18] is more effective on detail recovery compared to all-pixel based method PS-FCN [8], as highlighted in the yellow boxes (Observation 1).

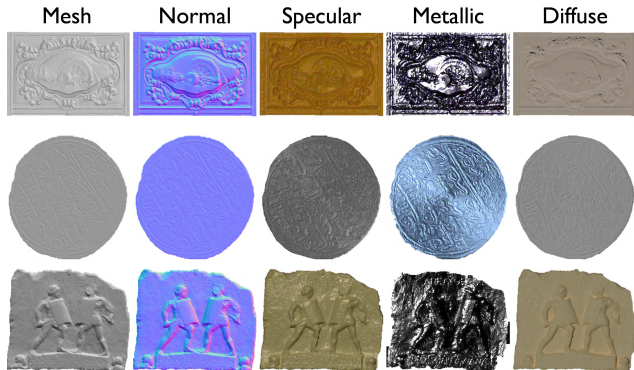


Figure 5: Overview of our synthetic dataset PS_RELIEF for fine-tuning photometric stereo on near-planar surfaces.

mean angular errors on DiLiGenT-II are significantly reduced after the fine-tuning on our near-planar dataset. However, we also observe performance degradation on DiLiGenT [41]. As shown in Fig. 6, the estimated surface normal for the BIRD in DiLiGenT-II is more reasonable after the fine-tuning, but with the consequence that the recovered CAT shape in DiLiGenT becomes more flattened. This behavior is not observed in calibrated photometric stereo methods such as PS-FCN, where the fine-tuning of PS-FCN

on PS_RELIEF does not influence the estimation on DiLiGenT [41] dramatically. We summarize the benchmark results of uncalibrated photometric stereo using DiLiGenT-II and fine-tuning existing methods using our synthetic PS_RELIEF dataset as the following observation:

Observation 2 *Near-planar surfaces are challenging for photometric stereo under uncalibrated light settings, and the normal estimation of learning-based uncalibrated photometric stereo is sensitive to the shape distributions present in the training dataset.*

4.3. Analysis to Different Reflectance Groups

As shown in Table 2, we find the four groups in DiLiGenT-II sorted by the MAnGE on different photometric stereo methods, arranged from high to low, are translucent, rough, specular, and metallic. The estimation errors come from non-Lambertian reflectance and surface geometry determining the shadows and inter-reflections. On the surface geometry side, we present the angular difference of a pure planar surface normal and the ground-truth surface normal in Table 2 (first column), showing that the shape variations of the translucent and the rough groups are greater than that of the metallic and specular groups. Also, surfaces in rough

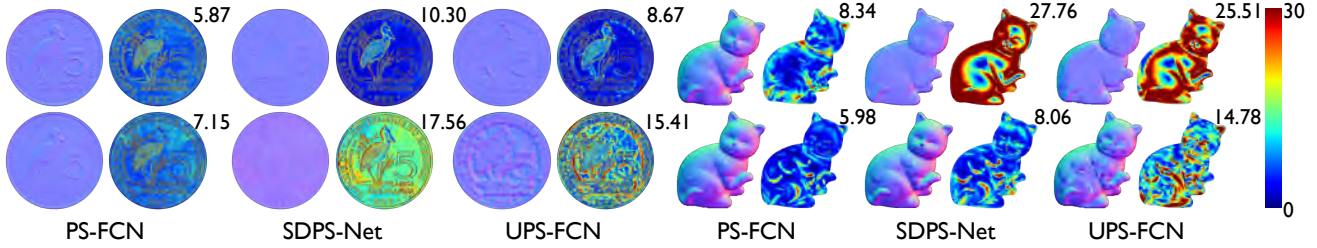


Figure 6: Surface normal estimates of photometric stereo methods with (top row) or without (bottom row) finetuning on PS_RELIEF. Compared to calibrated photometric stereo PS-FCN [8], learning-based uncalibrated photometric stereo (e.g. UPS-FCN [8], SDPS-Net [7]) are heavily influenced by the shape prior learned from the training dataset (Observation 2).

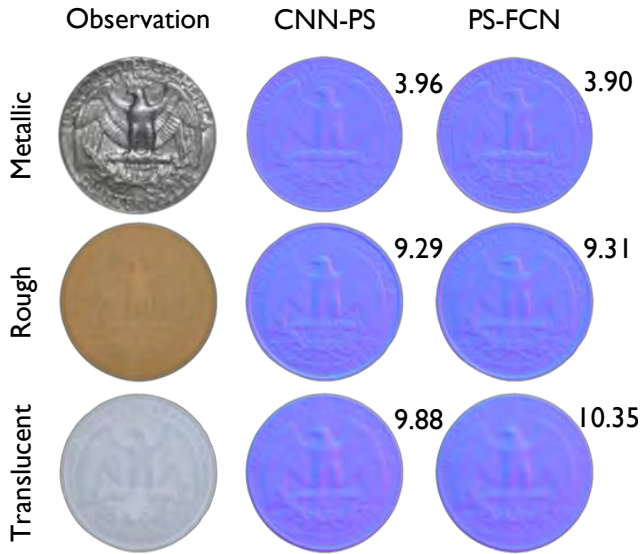


Figure 7: Ablation analysis to the influence of reflectance based on the same shape. The columns, from left to right, display the image observations, the estimated normal maps from CNN-PS [18], and from PS-FCN [8], with MAnGE labeled at the top-right corner. The rough and translucent materials are more challenging than the metallic materials for surface detail recovery (Observation 3).

and translucent groups contain more shadows and inter-reflections as their depth variation measured by PV is 4 mm compared to 1.5 mm, and 1 mm in the specular and metallic groups. On the reflectance side, the sub-surface scattering in the translucent group blurs the details of normal estimates as visualized in Fig. 4, resulting in greater estimation errors compared to the metallic and specular reflectance.

To disentangle the influence of reflectance and surface geometry on surface normal estimation, we conduct an ablation study on a metallic coin shown in Fig. 7. We scan this coin’s 3D mesh and create another two objects by 3D printing using translucent and rough materials that are the same as those used in our DiLiGenT-II. Based on the same

geometry, the surface normal estimation errors from CNN-PS [18] and PS-FCN [8] are higher on rough and translucent surfaces compared to metallic surfaces. Although the recovered details from the rough surface are sharper than the translucent one, the rough surface is still challenging due to its complex reflectance, and no previous photometric stereo method has targeted this kind of reflectance existing in objects covered by matte spray. We summarize these analysis results as our last observation:

Observation 3 *Near-planar surface normal estimation using photometric stereo methods remains a challenging task for translucent and rough surfaces, where surface details are significantly blurred due to the subsurface scattering in translucent surfaces.*

5. Conclusion

This paper builds a real-world photometric stereo dataset DiLiGenT-II focusing on near-planar surfaces with rich details, which are important to show the core strength of photometric stereo. We conduct benchmark evaluations on the dataset and draw three key observations. However, the evaluation metric utilized in this study is MAnGE, which assigns equal weights to surface normals regardless of the spatial distribution of surface details. Therefore, it is desired to devise a new evaluation metric that can measure the performance of surface detail recovery. Overall, we hope that DiLiGenT-II and the key observations will offer useful insights to further photometric stereo methods for detailed recovery of near-planar surfaces.

Acknowledgement

This work is supported by the National Key Research and Development Program of China under Grant No. 2019YFA0706701, National Natural Science Foundation of China under Grant No. 62136001, 62088102, 52175477. Heng Guo was supported by JSPS KAKENHI (Grant No. JP23H05491). The authors also want to thank *open-bayes.com* for providing part of the computing resource.

References

- [1] Jens Ackermann, Fabian Langguth, Simon Fuhrmann, and Michael Goesele. Photometric stereo for outdoor webcams. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012. 3
- [2] Neil Alldrin, Todd Zickler, and David Kriegman. Photometric stereo with non-parametric and spatially-varying reflectance. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2008. 3
- [3] Neil G. Alldrin, Satya P. Mallick, and David J. Kriegman. Resolving the generalized bas-relief ambiguity by entropy minimization. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2007. 4
- [4] Neil G. Alldrin, Todd Zickler, and David J. Kriegman. Photometric stereo with non-parametric and spatially-varying reflectance. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2008. 3
- [5] Brent Burley and Walt Disney Animation Studios. Physically-based shading at Disney. In *Proc. of SIGGRAPH*, 2012. 3, 7
- [6] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. ShapeNet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 3
- [7] Guanying Chen, Kai Han, Boxin Shi, Yasuyuki Matsushita, and Kwan-Yee K. Wong. Self-calibrating deep photometric stereo networks. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 4, 6, 7, 9
- [8] Guanying Chen, Kai Han, and Kwan-Yee K. Wong. PS-FCN: A flexible learning framework for photometric stereo. In *Proc. of European Conference on Computer Vision (ECCV)*, 2018. 3, 6, 7, 8, 9
- [9] Guanying Chen, Michael Waechter, Boxin Shi, Kwan-Yee K Wong, and Yasuyuki Matsushita. What is learned in deep uncalibrated photometric stereo? In *Proc. of European Conference on Computer Vision (ECCV)*, 2020. 4, 6, 7
- [10] Lixiong Chen, Yinqiang Zheng, Boxin Shi, Art Subpa-Asa, and Imari Sato. A microfacet-based reflectance model for photometric stereo with highly specular surfaces. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017. 3
- [11] Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2021. 2, 5, 7
- [12] Paul E Debevec and Jitendra Malik. Recovering high dynamic range radiance maps from photographs. In *ACM SIGGRAPH 2008 classes*. 2008. 5
- [13] Ondřej Drbohlav and Radim Šára. Specularities reduce ambiguity of uncalibrated photometric stereo. In *Proc. of European Conference on Computer Vision (ECCV)*, 2002. 4
- [14] Kenji Enomoto, Michael Waechter, Kiriakos N Kutulakos, and Yasuyuki Matsushita. Photometric stereo via discrete hypothesis-and-test search. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3
- [15] Dan B Goldman, Brian Curless, Aaron Hertzmann, and Steven M Seitz. Shape and spatially-varying brdfs from photometric stereo. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2009. 3
- [16] Bjoern Haefner, Songyou Peng, Alok Verma, Yvain Quéau, and Daniel Cremers. Photometric depth super-resolution. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. 3
- [17] Yannick Hold-Geoffroy, Paulo Gotardo, and Jean-François Lalonde. Single day outdoor photometric stereo. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. 3
- [18] Satoshi Ikehata. CNN-PS: CNN-based photometric stereo for general non-convex surfaces. In *Proc. of European Conference on Computer Vision (ECCV)*, 2018. 3, 6, 7, 8, 9
- [19] Satoshi Ikehata. PS-Transformer: Learning sparse photometric stereo network using self-attention mechanism. *Proc. of the British Machine Vision Conference (BMVC)*, 2022. 4
- [20] Satoshi Ikehata. Scalable, detailed and mask-free universal photometric stereo. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 4, 6
- [21] Satoshi Ikehata, David Wipf, Yasuyuki Matsushita, and Kiyoharu Aizawa. Robust photometric stereo using sparse regression. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012. 7
- [22] Wenzel Jakob. Mitsuba renderer, 2010. 2
- [23] Micah K Johnson and Edward H Adelson. Shape estimation in natural illumination. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2011. 3
- [24] Yakun Ju, Kin-Man Lam, Yang Chen, Lin Qi, and Junyu Dong. Pay attention to devils: A photometric stereo network for better details. In *Proc. of International Joint Conference on Artificial Intelligence*, 2021. 4
- [25] Yakun Ju, Kin-Man Lam, Wuyuan Xie, Huiyu Zhou, Junyu Dong, and Boxin Shi. Deep learning methods for calibrated photometric stereo and beyond: A survey. *arXiv preprint arXiv:2212.08414*, 2022. 3, 6
- [26] Yakun Ju, Boxin Shi, Muwei Jian, Lin Qi, Junyu Dong, and Kin-Man Lam. NormAttention-PSN: A high-frequency region enhanced photometric stereo network with normalized attention. *International Journal of Computer Vision*, 2022. 4, 6, 7
- [27] Berk Kaya, Suryansh Kumar, Carlos Oliveira, Vittorio Ferrari, and Luc Van Gool. Uncalibrated neural inverse rendering for photometric stereo of general surfaces. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3, 4
- [28] Junxuan Li, Antonio Robles-Kelly, Shaodi You, and Yasuyuki Matsushita. Learning to minify photometric stereo. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 3
- [29] Min Li, Zhenglong Zhou, Zhe Wu, Boxin Shi, Changyu Diao, and Ping Tan. Multi-view photometric stereo: A robust solution and benchmark dataset for spatially varying isotropic materials. *IEEE Transactions on Image Processing*, 2020. 1
- [30] Min Li, Zhenglong Zhou, Zhe Wu, Boxin Shi, Changyu Diao, and Ping Tan. Multi-view photometric stereo: A robust solution and benchmark dataset for spatially varying

- isotropic materials. *IEEE Transactions on Image Processing*, 2020. 3
- [31] Daniel Lichy, Jiaye Wu, Soumyadip Sengupta, and David W Jacobs. Shape and material capture at home. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 4, 6
- [32] Fotios Logothetis, Ignas Budvytis, Roberto Mecca, and Roberto Cipolla. PX-NET: simple and efficient pixel-wise training of photometric stereo networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021. 3, 6
- [33] Wojciech Matusik, Hanspeter Pfister, Matt Brand, and Leonard McMillan. A data-driven reflectance model. In *Proc. of SIGGRAPH*, 2003. 3
- [34] Roberto Mecca, Fotios Logothetis, Ignas Budvytis, and Roberto Cipolla. LUCES: A dataset for near-field point light source photometric stereo. In *Proc. of the British Machine Vision Conference (BMVC)*, 2021. 1, 3, 5
- [35] Michael Oren and Shree K Nayar. Generalization of Lambert’s reflectance model. In *Proc. of Annual conference on Computer Graphics and Interactive Techniques*, 1994. 1, 4
- [36] Thoma Papadhimetri and Paolo Favaro. A new perspective on uncalibrated photometric stereo. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013. 4
- [37] Thoma Papadhimetri and Paolo Favaro. A closed-form, consistent and robust solution to uncalibrated photometric stereo via local diffuse reflectance maxima. *International Journal of Computer Vision*, 2014. 4, 6, 7
- [38] Yvain Quéau, Tao Wu, François Lauze, Jean-Denis Durou, and Daniel Cremers. A non-convex variational approach to photometric stereo under inaccurate lighting. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 3
- [39] Jiejie Ren, Feishi Wang, Jiahao Zhang, Qian Zheng, Mingjun Ren, and Boxin Shi. DiLiGenT10²: A photometric stereo benchmark dataset with controlled shape and material variation. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1, 3, 5, 7
- [40] Hiroaki Santo, Masaki Samejima, Yusuke Sugano, Boxin Shi, and Yasuyuki Matsushita. Deep photometric stereo network. In *Proc. of International Conference on Computer Vision Workshops (ICCVW)*, 2017. 3
- [41] 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 Transactions on Pattern Analysis and Machine Intelligence*, 2019. 1, 3, 5, 6, 7, 8
- [42] Boxin Shi, Ping Tan, Yasuyuki Matsushita, and Katsushi Ikeuchi. Bi-polynomial modeling of low-frequency reflectances. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2014. 3, 5, 6, 7
- [43] Boxin Shi, Zhe Wu, Zhipeng Mo, Dinglong Duan, Sai-Kit Yeung, and Ping Tan. A benchmark dataset and evaluation for non-Lambertian and uncalibrated photometric stereo. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1
- [44] William M. Silver. *Determining shape and reflectance using multiple images*. PhD thesis, Massachusetts Institute of Technology, 1980. 1
- [45] Bo Sun, Kalyan Sunkavalli, Ravi Ramamoorthi, Peter N Belhumeur, and Shree K Nayar. Time-varying brdfs. *IEEE Transactions on Visualization and Computer Graphics*, 2007. 4
- [46] Kenneth E Torrance and Ephraim M Sparrow. Theory for off-specular reflection from roughened surfaces. *Journal of the Optical Society of America*, 1967. 4
- [47] Robert J. Woodham. Photometric method for determining surface orientation from multiple images. *Optical engineering*, 1980. 1, 5, 6, 7
- [48] Lun Wu, Arvind Ganesh, Boxin Shi, Yasuyuki Matsushita, Yongtian Wang, and Yi Ma. Robust photometric stereo via low-rank matrix completion and recovery. In *Proc. of Asian Conference on Computer Vision (ACCV)*, 2010. 3
- [49] Tai-Pang Wu and Chi-Keung Tang. Photometric stereo via expectation maximization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2010. 5, 6, 7
- [50] Ying Xiong, Ayan Chakrabarti, Ronen Basri, Steven J Gortler, David W Jacobs, and Todd Zickler. From shading to local shape. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2014. 3
- [51] Ying Xiong, Ayan Chakrabarti, Ronen Basri, Steven J. Gortler, David W. Jacobs, and Todd Zickler. From shading to local shape. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2015. 3
- [52] Zhuokun Yao, Kun Li, Ying Fu, Haofeng Hu, and Boxin Shi. GPS-Net: Graph-based photometric stereo network. *Advances in Neural Information Processing Systems*, 2020. 4, 6
- [53] Qian Zheng, Yiming Jia, Boxin Shi, Xudong Jiang, Ling-Yu Duan, and Alex C. Kot. SPLINE-Net: Sparse photometric stereo through lighting interpolation and normal estimation networks. In *Proc. of International Conference on Computer Vision (ICCV)*, 2019. 3
- [54] Qian Zheng, Boxin Shi, and Gang Pan. Summary study of data-driven photometric stereo methods. *Virtual Reality & Intelligent Hardware*, 2020. 2, 3, 6