

000
001
002003

DiLiGenT-II: Photometric Stereo for Planar Surfaces with Rich Details – 004 Benchmark Dataset and Beyond

005
006
007008 Anonymous ICCV submission
009010
011012 Paper ID 1565
013014
015

Abstract

016
017

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 <http://be.available.upon.acceptance>.

034
035

1. Introduction

036
037

Photometric stereo [45, 42] 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.

045
046
047
048
049
050
051
052

To evaluate the effectiveness of photometric stereo methods, a batch of real-world benchmark datasets have been built such as DiLiGenT [41, 39] (and its multi-view extension DiLiGenT-MV [28]), DiLiGenT10² [37] for distant lights, and LUCES [32] 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 [33]. In addition, our DiLiGenT-II takes an optical profilometer to capture the ultra-precise surface 3D structure in nanometer 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 [52]; the challenging problems of uncalibrated photometric stereo for near-planar surfaces; and the influence of translucent and rough reflectance on surface detail recovery. The analy-

¹“DiLiGenT” [39] 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.

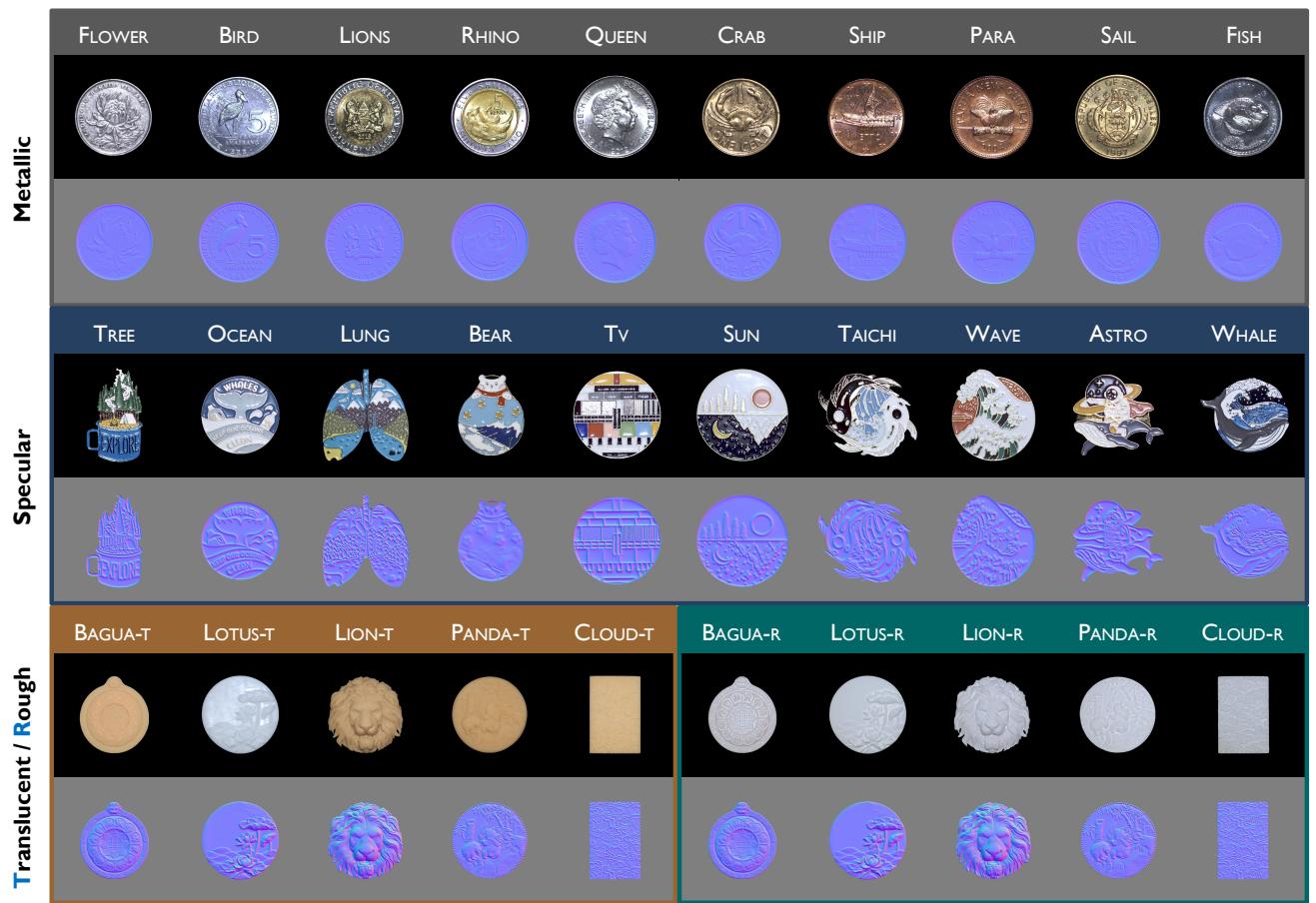


Figure 1: Overview of DiLiGenT-II. We collect 4 groups of near-planar objects with rich surface details and diverse reflectance types (metallic, specular, rough, and translucent). The corresponding ‘ground-truth’ surface normals shown in the even rows are measured via a precise profilometer in the accuracy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

sis 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 [21] and Blender [11] to create image observations and the corresponding surface normal maps of diverse synthetic scenes. DPSN [38] introduces the first synthetic photometric stereo dataset BlobbyPS, comprising 10 smooth Blobby [22] shapes with reflectance assigned by measured MERL BRDFs [31] 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 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 [22] and Sculpture dataset [6], leading to a shape domain gap to flatten surfaces. In this paper, we render a synthetic dataset contain-

216 Table 1: Summary of real-word photometric stereo datasets. Material: controlled (C: fabricated or carefully selected with
 217 controlled categories) or uncontrolled (UC: randomly picked up from daily objects); Ground Truth (GT measure): from
 218 CAD/Scanned models with registration (+Reg) or from photometric stereo (PS). Number (#) of shapes, lights, and sets (one
 219 set means a sequence of photometric stereo images under varying lighting conditions used for computation).
 220

Dataset							
DiLiGenT-II	DiLiGenT10 ² [37]	DiLiGenT [39]	LUCES [32]	Harvard [49]	ETHz [26]	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

232 ing 127 near-planar surfaces with rich details to enhance
 233 learning-based photometric stereo on near-planar surface
 234 recovery.

235 **Real-world datasets** complement the gap between computer
 236 graphics rendering and real-world imaging process.
 237 The Gourd&Apple dataset [2] releases image observations
 238 of 2 objects with spatially-varying isotropic BRDFs. Harvard
 239 dataset [48] contains 7 surfaces with uniform diffuse
 240 reflectance. However, the ground truth surface normal of
 241 the above two datasets is not provided. DiLiGenT [39]
 242 records 10 objects with diverse shapes and general non-
 243 Lambertian materials. Starting from DiLiGenT [39], bench-
 244 mark evaluation of photometric stereo becomes available
 245 based on the ‘ground truth’ surface normal from scanned
 246 meshes. The following-up datasets further extends DiLi-
 247 GenT [39] from the perspective of multi-views (DiLiGenT-MV
 248 [29]), controlled materials and shapes (DiLiGenT
 249 10² [37]), near-field illumination (LUCES [32]), environ-
 250 ment illumination [1, 16, 17], and global illumination ef-
 251 fects [26].

252 As summarized in Tab. 1, objects contained in existing
 253 real-world photometric stereo datasets are mostly bulgy
 254 and smooth, which do not include flattened or near-planar
 255 surfaces though they are commonly seen in our daily life.
 256 More importantly, the smooth target shapes are not suitable
 257 to evaluate the uniqueness of photometric stereo over other
 258 3D reconstruction techniques, that is, the high-fidelity sur-
 259 face detail recovery, especially for delicate structures. To
 260 address these two problems, we newly build a real-world
 261 photometric stereo dataset containing near-planar surfaces
 262 with rich geometric details.

263 2.2. Photometric Stereo Methods

264 We briefly review existing photometric stereo meth-
 265 ods based on distant calibrated and uncalibrated light set-
 266 tings. Please refer to the photometric stereo survey of

287 non-learning based methods [39] and learning-based meth-
 288 ods [52, 24] for more comprehensive analysis.

289 **Calibrated** photometric stereo works with calibrated dis-
 290 tant lights. Existing non-learning based methods either treat
 291 specular highlights and shadows as sparse outliers [46, 36]
 292 or propose parametric or data-based reflectance model to
 293 explicitly handle the non-Lambertian reflection [10, 14,
 294 40, 15]. Beginning from DPSN [38], the image observa-
 295 tions are directly mapped to the corresponding surface nor-
 296 mal via deep neural networks, where the non-Lambertian
 297 reflectance is implicitly learned from the synthetic train-
 298 ing data. The network structure in existing learning-based
 299 methods can be divided into the all-pixel branch (represen-
 300 tative work: PS-FCN [8]) and the per-pixel branch (represen-
 301 tative work: CNN-PS [18]). Based on these two typ-
 302 ical network structures, following-up works further pro-
 303 mote photometric stereo by addressing the global illumi-
 304 nation effects (*e.g.* PX-Net [30]), reducing the number of
 305 inputs (*e.g.* SPLINE-Net [51], LMPS [27]), combining the
 306 merit of per-pixel and all-pixel methods (*e.g.* GPS-Net [50],
 307 PS-Transformer [19]). However, few methods focus on the
 308 recovery of surface details despite its importance in pho-
 309 tomatic stereo. Ju *et al.* [23] firstly addressed the high-
 310 frequency surface details in photometric stereo based on
 311 attention-weighted loss. Their following-up work [25] fur-
 312 ther enhanced the detail recovery accuracy via a double-gate
 313 normalization and a parallel high-resolution structure.

314 **Uncalibrated** photometric stereo works under unknown
 315 distant light directions, so that photometric stereo becomes
 316 more challenging even with Lamebrain reflectance assump-
 317 tion. Non-learning based methods adopt the local diffuse
 318 reflectance (LDR) maxima [35], perspective camera [34],
 319 specularities [13], albedo entropy [3] to resolve the sur-
 320 face normal estimation ambiguity in uncalibrated Lamber-
 321 tian photometric stereo. Beginning from SDPS-Net [7],
 322 learning-based uncalibrated photometric stereo methods ca-

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

pable 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.* [26] 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.

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.

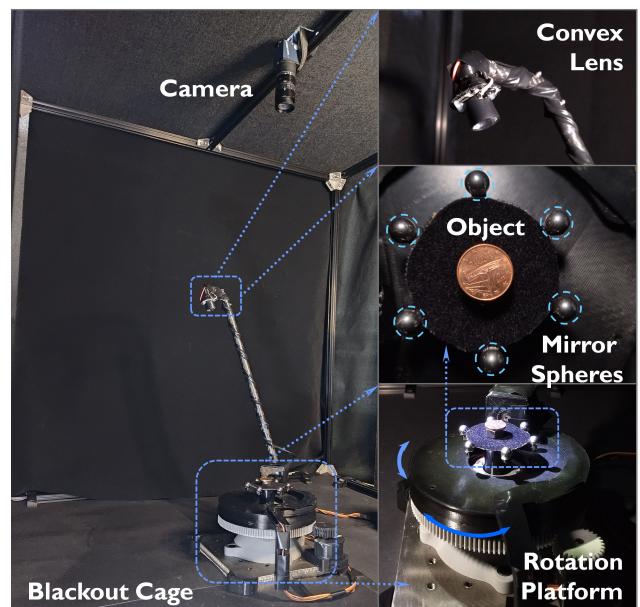
3.1. Objects groups

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 $15 \sim 25$ 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,



378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

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’.

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 [33] and the Torrance-Sparrow specular model [44], as analyzed in [43].

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 sin-

²Kexcelled Ultradetail: <https://www.kexcelled3d.com/products/ultrdetial/>. Retrieved March 6, 2023.

³FA-5: <http://cysygroup.com/en/product.asp?category=NDT&page=5>. Retrieved March 6, 2023.



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

gle 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 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 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 [39, 37]. Please check our supplementary material for more details.

3.3. Obtaining ‘GT’ Normal Map

Similar to previous real-world datasets: DiLiGenT [39] and LUCES [32], 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

⁴Camera website: <https://www.daheng-imaging.com/product/area-scan-cameras/daheng/mer-u3/2592.html>. Retrieved March 6, 2023.

⁵Alicona website: <https://www.alicona.com/en/product/infinitefocus>. Retrieved March 6, 2023.

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 [39, 37].

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 [39] adopt non-Learning photometric stereo, we choose the baseline method (least-square based Lambertian photometric stereo [45], LSPS), baseline method with position thresholding strategy [39] (TH28 and TH46 reject pixels whose intensities under varying lights are outside the range of [20%, 80%] and [40%, 60%]), respectively), WG10 [47] (a robust photometric stereo method based on outlier rejection), ST14 [40] (showing best performance at DiLiGenT reported in [39]). Based on surveys about the learning-based photometric stereo [52, 24], we choose representative networks handling photometric stereo in an all-pixel manner (NormAttention-PSN [25], PS-FCN [8]), per-pixel manner (CNN-PS [18], PX-Net [30]), and the method GPS-Net [50] that considers both the per-pixel and the all-pixel structures. Besides the above calibrated photometric stereo methods, we also evaluate existing uncalibrated photometric stereo approaches, including a non-learning based method PF14 [35] (showing the best performance under uncalibrated setting as reported in [39]), three representative learning-based methods SDPS-Net [7], UPS-FCN [8] and UPS-GCNet [9] published in recent years. During the evaluation, we adopted the code and pre-trained model released by the authors to process the collected data in DiLiGenT-II.

Similar to DiLiGenT [39] and DiLiGenT10² [37], 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.

4.2. Analysis to Different Baselines

As shown in Fig. 4, we evaluate the surface normal estimation accuracy for all methods on DiLiGenT-II. In the following, we first analyze the surface detail recovery from

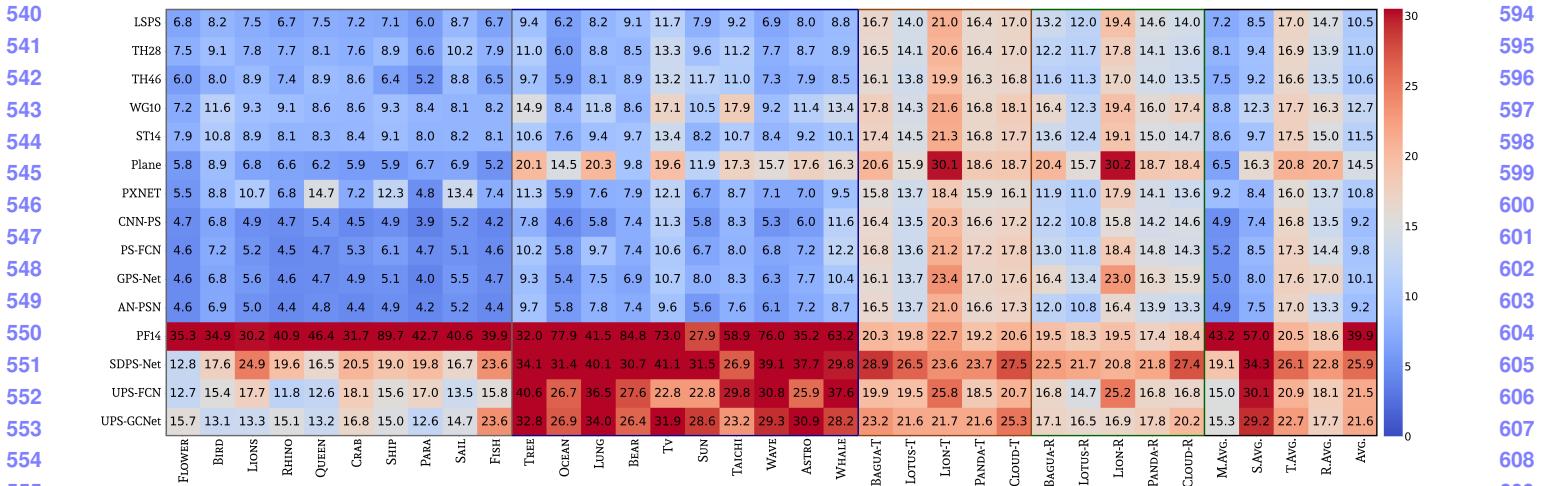


Figure 4: Benchmark results on our real-world dataset DiLiGenT-II. Mean angular error in degree of each object on various methods are presented, and the average angular error of each material group is concluded in the rightmost columns (best zoom in and viewed in color). We denote ‘NA-PSN’ as the abbreviation of NormAttention-PSN [25].

calibrated photometric stereo, followed by the analysis of uncalibrated photometric stereo on near-planar surfaces.

Calibrated photometric stero. From Fig. 4, we observe that the LSPS [45] shows the best performance over other non-learning-based photometric stereo methods, while the angular error difference between LSPS [45], TH28, TH46, and ST14 [40] 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 [45]. 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 [37].

Among learning-based calibrated photometric stereo, CNN-PS [18] and NormAttention-PSN [25] achieve smaller mean angular errors, showing their advantages on detailed surface recovery. As shown in Fig. 5, we visualize the error distributions of PS-FCN [8], CNN-PS [20], and NormAttention-PSN [25] on four representative surfaces belong to the four reflectance groups. PS-FCN [7] outputs blurry normal estimation results as highlighted in Fig. 5, 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 [25] 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 stero. From Fig. 4, 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 [35], 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. 6, 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

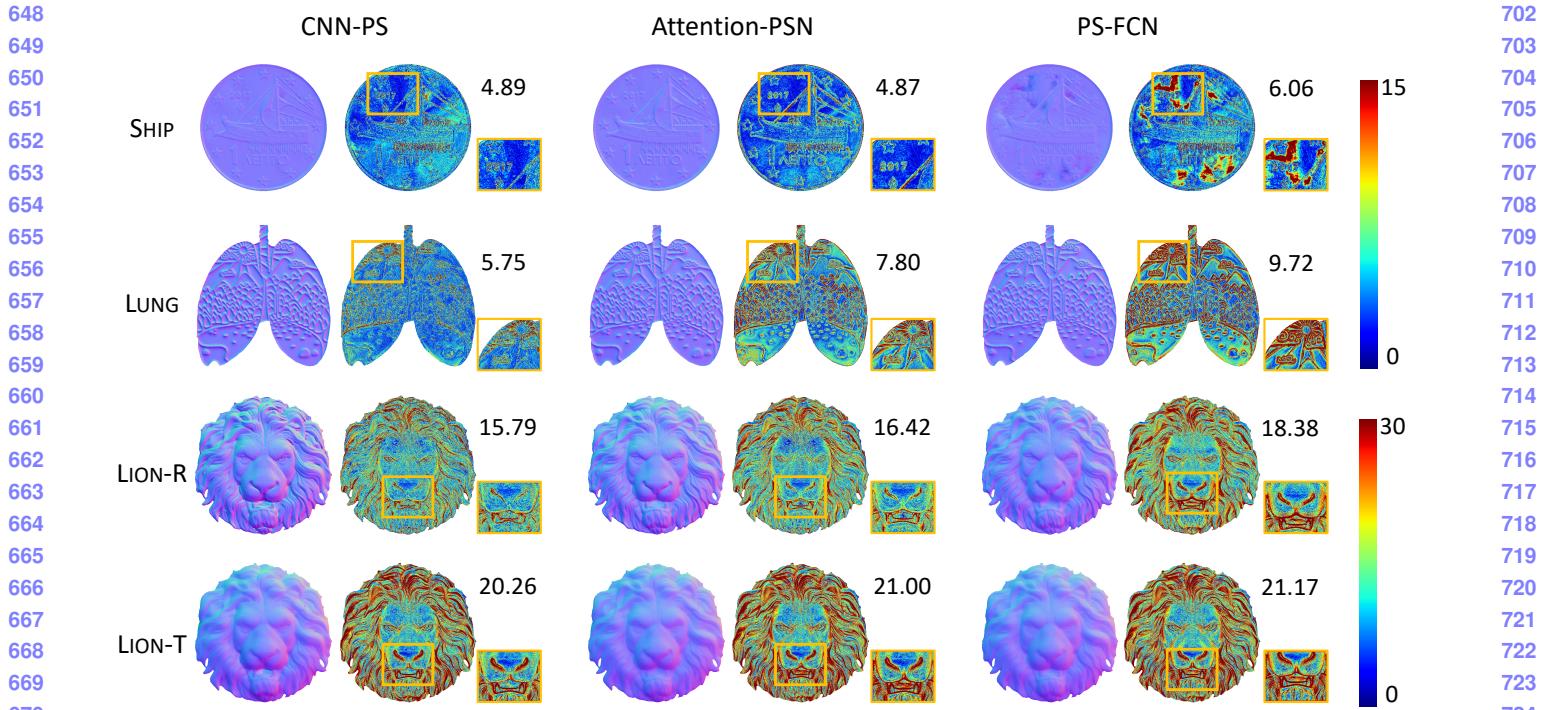


Figure 5: Visualization of surface detail recovery from different photometric stereo methods, where the odd and even columns plot the estimated surface normals and the corresponding angular error maps, respectively. Per-pixel based photometric stereo method CNN-PS [18] is more effective on detail recovery compared to all-pixel based method PS-FCN [8], as highlighted in the yellow boxes (Observation 1).

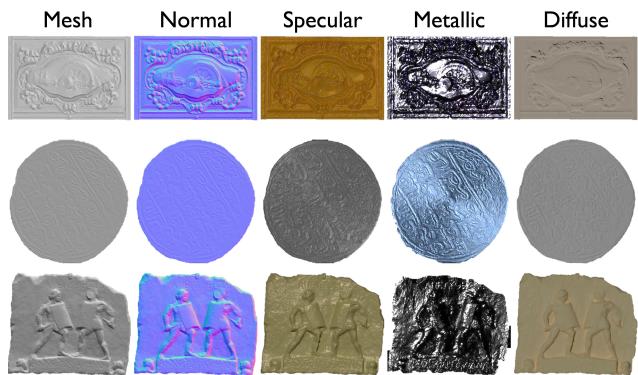


Figure 6: Overview of our synthetic dataset PS_RELIEF for fine-tuning photometric stereo on near-planar surfaces.

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 DiLiGenT^{10²} [37].

As shown in Table 2, we fine-tune the UPS-FCN [8] and SDPS-Net [7] with PS_RELIEF and test it on real-world dataset DiLiGenT [39] and DiLiGenT-II, which contain bulgy shapes and near-planar shapes, respectively. The

Table 2: Ablation study on photometric stereo methods trained with or without fine-tuning (FT.) on the near-planar synthetic dataset PS_RELIEF.

Light	Method	DiLiGenT [39]		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

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 [39]. As shown in Fig. 7, 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 [39] 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:

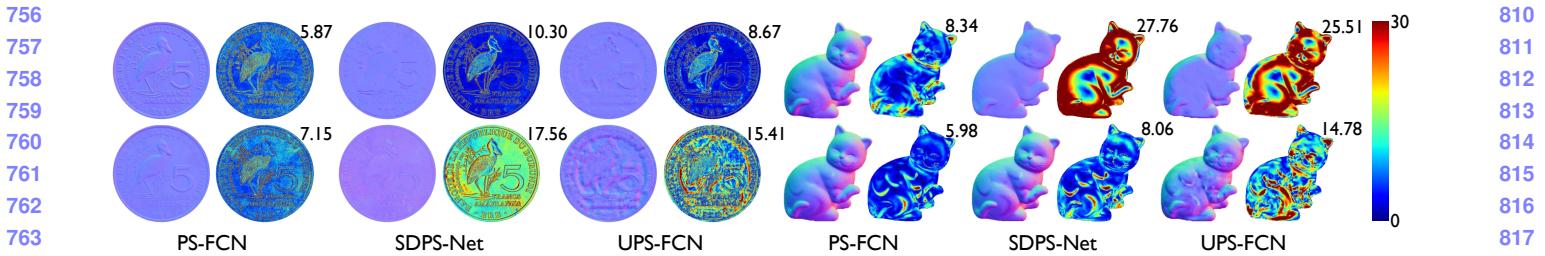


Figure 7: Surface normal estimates of photometric stereo methods with (top row) or without (bottom row) finetuning on PS_RELIEF. Compared to calibrated photometric stereo PS-FCN [8], learning-based uncalibrated photometric stereo (e.g. UPS-FCN [8], SDPS-Net [7]) are heavily influenced by the shape prior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

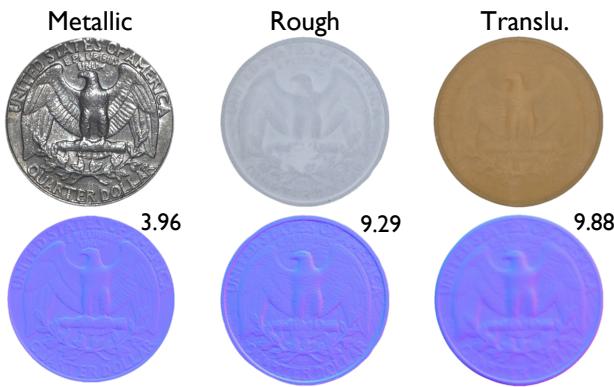


Figure 8: Ablation analysis to the influence of reflectance based on the same shape. The top and bottom rows show the image observations and the estimated normal maps from CNN-PS [18], with MAngE labeled at the top-right corner. The rough and translucent materials are still challenging for recovering surface details (Observation 3).

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 Fig. 4, 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 Fig. 4 (6-th row), 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 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. 5, 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. 8. 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] 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.

864

References

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

- [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. 3
- [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. 2, 6
- [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. 2
- [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. 3, 4, 5, 6, 7, 8
- [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. 2, 3, 5, 6, 7, 8
- [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, 5, 6
- [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. 3
- [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. 2, 3, 5, 6, 7, 8
- [19] Satoshi Ikehata. PS-Transformer: Learning sparse photometric stereo network using self-attention mechanism. *Proc. of the British Machine Vision Conference (BMVC)*, 2022. 3
- [20] 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. 6
- [21] Wenzel Jakob. Mitsuba renderer, 2010. 2
- [22] 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. 2
- [23] 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. 3
- [24] 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, 5
- [25] 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. 3, 5, 6
- [26] 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
- [27] 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
- [28] 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
- [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

- 972 isotropic materials. *IEEE Transactions on Image Processing*, 2020. 3
- 973 [30] Fotios Logothetis, Ignas Budvytis, Roberto Mecca, and
974 Roberto Cipolla. PX-NET: simple and efficient pixel-wise
975 training of photometric stereo networks. In *Proceedings of
976 the IEEE/CVF International Conference on Computer Vision*, 2021. 3, 5
- 977 [31] Wojciech Matusik, Hanspeter Pfister, Matt Brand, and
978 Leonard McMillan. A data-driven reflectance model. In
979 *Proc. of SIGGRAPH*, 2003. 2
- 980 [32] Roberto Mecca, Fotios Logothetis, Ignas Budvytis, and
981 Roberto Cipolla. LUCES: A dataset for near-field point light
982 source photometric stereo. In *Proc. of the British Machine
983 Vision Conference (BMVC)*, 2021. 1, 3, 5
- 984 [33] Michael Oren and Shree K Nayar. Generalization of Lam-
985 bert's reflectance model. In *Proc. of Annual conference on
986 Computer Graphics and Interactive Techniques*, 1994. 1, 4
- 987 [34] Thoma Papadimitri and Paolo Favaro. A new perspective
988 on uncalibrated photometric stereo. In *Proc. of IEEE Confer-
989 ence on Computer Vision and Pattern Recognition (CVPR)*,
990 2013. 3
- 991 [35] Thoma Papadimitri and Paolo Favaro. A closed-form, con-
992 sistent and robust solution to uncalibrated photometric stereo
993 via local diffuse reflectance maxima. *International Journal
994 of Computer Vision*, 2014. 3, 5, 6
- 995 [36] Yvain Quéau, Tao Wu, François Lauze, Jean-Denis Durou,
996 and Daniel Cremers. A non-convex variational approach to
997 photometric stereo under inaccurate lighting. In *Proc. of
998 IEEE Conference on Computer Vision and Pattern Recog-
999 nition (CVPR)*, 2017. 3
- 1000 [37] Jieji Ren, Feishi Wang, Jiahao Zhang, Qian Zheng, Mingjun
1001 Ren, and Boxin Shi. DiLiGenT10²: A photometric stereo
1002 benchmark dataset with controlled shape and material varia-
1003 tion. In *Proc. of IEEE Conference on Computer Vision and
1004 Pattern Recognition (CVPR)*, 2022. 1, 3, 5, 6, 7
- 1005 [38] Hiroaki Santo, Masaki Samejima, Yusuke Sugano, Boxin
1006 Shi, and Yasuyuki Matsushita. Deep photometric stereo net-
1007 work. In *Proc. of International Conference on Computer
1008 Vision Workshops (ICCVW)*, 2017. 2, 3
- 1009 [39] Boxin Shi, Zhipeng Mo, Zhe Wu, Dinglong Duan, Sai-Kit
1010 Yeung, and Ping Tan. A benchmark dataset and evaluation
1011 for non-Lambertian and uncalibrated photometric stereo.
1012 *IEEE Transactions on Pattern Analysis and Machine Intel-
1013 ligence*, 2019. 1, 3, 5, 7
- 1014 [40] Boxin Shi, Ping Tan, Yasuyuki Matsushita, and Katsushi
1015 Ikeuchi. Bi-polynomial modeling of low-frequency re-
1016 flectances. *IEEE Transactions on Pattern Analysis and Ma-
1017 chine Intelligence*, 2014. 3, 5, 6
- 1018 [41] Boxin Shi, Zhe Wu, Zhipeng Mo, Dinglong Duan, Sai-Kit
1019 Yeung, and Ping Tan. A benchmark dataset and evaluation
1020 for non-Lambertian and uncalibrated photometric stereo. In
1021 *Proc. of IEEE Conference on Computer Vision and Pattern
1022 Recognition (CVPR)*, 2016. 1
- 1023 [42] William M. Silver. *Determining shape and reflectance us-
1024 ing multiple images*. PhD thesis, Massachusetts Institute of
Technology, 1980. 1
- 1025 [43] Bo Sun, Kalyan Sunkavalli, Ravi Ramamoorthi, Peter N
Belhumeur, and Shree K Nayar. Time-varying brdfs.
- 1026 *IEEE Transactions on Visualization and Computer Graph-
1027 ics*, 2007. 4
- 1028 [44] Kenneth E Torrance and Ephraim M Sparrow. Theory for
1029 off-specular reflection from roughened surfaces. *Journal of
1030 the Optical Society of America*, 1967. 4
- 1031 [45] Robert J. Woodham. Photometric method for determining
1032 surface orientation from multiple images. *Optical engineer-
1033 ing*, 1980. 1, 5, 6
- 1034 [46] Lun Wu, Arvind Ganesh, Boxin Shi, Yasuyuki Matsushita,
1035 Yongtian Wang, and Yi Ma. Robust photometric stereo via
1036 low-rank matrix completion and recovery. In *Proc. of Asian
1037 Conference on Computer Vision (ACCV)*, 2010. 3
- 1038 [47] Tai-Pang Wu and Chi-Keung Tang. Photometric stereo via
1039 expectation maximization. *IEEE Transactions on Pattern
1040 Analysis and Machine Intelligence*, 2010. 5
- 1041 [48] Ying Xiong, Ayan Chakrabarti, Ronen Basri, Steven J
1042 Gortler, David W Jacobs, and Todd Zickler. From shading
1043 to local shape. *IEEE Transactions on Pattern Analysis and
1044 Machine Intelligence*, 2014. 3
- 1045 [49] Ying Xiong, Ayan Chakrabarti, Ronen Basri, Steven J.
1046 Gortler, David W. Jacobs, and Todd Zickler. From shading
1047 to local shape. *IEEE Transactions on Pattern Analysis and
1048 Machine Intelligence*, 2015. 3
- 1049 [50] Zhuokun Yao, Kun Li, Ying Fu, Haofeng Hu, and Boxin
1050 Shi. GPS-Net: Graph-based photometric stereo network.
1051 *Advances in Neural Information Processing Systems*, 2020. 3,
5
- 1052 [51] Qian Zheng, Yiming Jia, Boxin Shi, Xudong Jiang, Ling-Yu
1053 Duan, and Alex C. Kot. SPLINE-Net: Sparse photometric
1054 stereo through lighting interpolation and normal estimation
1055 networks. In *Proc. of International Conference on Computer
1056 Vision (ICCV)*, 2019. 3
- 1057 [52] Qian Zheng, Boxin Shi, and Gang Pan. Summary study of
1058 data-driven photometric stereo methods. *Virtual Reality &
1059 Intelligent Hardware*, 2020. 1, 3, 5