Single Day Outdoor Photometric Stereo

Yannick Hold-Geoffroy[®], Paulo Gotardo, and Jean-François Lalonde

Abstract—Photometric Stereo (PS) under outdoor illumination remains a challenging, ill-posed problem due to insufficient variability in illumination. Months-long capture sessions are typically used in this setup, with little success on shorter, single-day time intervals. In this paper, we investigate the solution of outdoor PS over a single day, under different weather conditions. First, we investigate the relationship between weather and surface reconstructability in order to understand when natural lighting allows existing PS algorithms to work. Our analysis reveals that partially cloudy days improve the conditioning of the outdoor PS problem while sunny days do not allow the unambiguous recovery of surface normals from photometric cues alone. We demonstrate that calibrated PS algorithms can thus be employed to reconstruct Lambertian surfaces accurately under partially cloudy days. Second, we solve the ambiguity arising in clear days by combining photometric cues with prior knowledge on material properties, local surface geometry and the natural variations in outdoor lighting through a CNN-based, weakly-calibrated PS technique. Given a sequence of outdoor images captured during a single sunny day, our method robustly estimates the scene surface normals with unprecedented quality for the considered scenario. Our approach does not require precise geolocation and significantly outperforms several state-of-the-art methods on images with real lighting, showing that our CNN can combine efficiently learned priors and photometric cues available during a single sunny day.

Index Terms—Photometric stereo, high dynamic range, deep learning, outdoor lighting

17 **1** INTRODUCTION

3

5

6

7

8 9

10 11

12 13

14

15

16

CINCE its inception in the early 80s, Photometric Stereo 18 $\mathcal{J}(PS)$ [1] has been explored under many an angle. 19 20 Whether it has been to improve its ability to deal with complex materials [2] or lighting conditions [3], the myriad of 21 22 papers published on the topic are testament to the interest this technique has garnered in the community. While most 23 of the papers on this topic have focused on images captured 24 in the lab, recent progress has allowed the application of PS 25 on images captured outdoors, lit by the more challenging 26 case of uncontrollable, natural illumination. 27

While capturing more data in the lab can be done rela-28 tively easily, the same cannot be said for outdoor imagery. 29 Indeed, one does not control the sun and the other atmo-30 31 spheric elements in the sky; so one must wait for lighting conditions to change on their own. A creative solution to 32 33 this problem was proposed in [4], but it is limited to objects that can be placed on a small moving platform. Therefore, 34 35 capturing more data for fixed, large objects still means waiting days, or even months, potentially [5], [6]. 36

A promising approach to answer this question is to use more elaborate models of illumination—high dynamic range (HDR) environment maps [7]—as input to outdoor PS. Favorable results have been reported in [8] for outdoor images taken in a single day, within an interval of just eight

- Y. Hold-Geoffroy is with the Adobe, San Jose, CA 95110-2704. E-mail: yannickhold@gmail.com.
- P. Gotardo is with the Disney Research Studios, 8006 Zurich, Switzerland. E-mail: paulo.gotardo@disneyresearch.com.
- J.-F. Lalonde is with the Université Laval, Quebec city, QC G1V 0A6, Canada. E-mail: jflalonde@gel.ulaval.ca.

Manuscript received 8 Mar. 2019; revised 21 Sept. 2019; accepted 11 Dec. 2019. Date of publication 0 . 0000; date of current version 0 . 0000. (Corresponding author: Yannick Hold-Geoffroy.) Recommended for acceptance by K. Nishino. Digital Object Identifier no. 10.1109/TPAMI.2019.2962693 hours. However, the quality of outdoor results is reported 42 to be inferior to that obtained in indoor environments. The 43 decline is attributed to modest variation in sunlight, but no 44 clear explanation is found in the literature. This observation 45 leads to many interesting, unanswered questions: had the 46 atmospheric conditions been different on that day, could 47 the quality of their results have been better? Is a full day of 48 observations enough to obtain good results in outdoor PS? 49

This paper investigates PS in outdoor environments over 50 the course of a single day and under a variety of sunlight 51 conditions. Our first goal is to assess the *reconstructibility* of 52 surface patches as a function of their orientation and the 53 illumination conditions. This is done using a large database 54 of sky probes [9], capturing the variability of natural, out-55 door illumination. A detailed look at the conditions under 56 which normals can be reconstructed reliably is presented, 57 followed by an analysis of surface reconstruction stability. 58

Our analysis reveals that reconstruction performance of 59 classical PS methods can be categorized in two different sky 60 types: partially cloudy and clear days. Interestingly, partially 61 cloudy days typically offer better reconstruction accuracy, 62 while clear days generally yield poor performance. During 63 clear days, photometric cues alone do not provide a stable 64 solution to the PS problem, leaving it under-constrained [1]. 65 Our insight to solve this issue is to augment the photometric 66 cues with learned features on geometry, reflectance and 67 lighting to resolve ambiguities in singe-day outdoor PS. 68

We summarize our contributions as follows:

 an analysis of the conditioning of outdoor PS given 70 photometric cues captured over a single day; 71

69

 a framework for predicting the performance of sin- 72 gle-day outdoor PS with calibrated lighting, and its 73 application in reconstructing surfaces on partially 74 cloudy days; 75

^{0162-8828 © 2019} IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information.



Fig. 1. Examples from our dataset of HDR outdoor illumination conditions. In all, our dataset contains 3,800 different illumination conditions, captured from 10:30 until 16:30, during 23 days, spread over ten months and at two geographical locations. Each image is stored in the 32-bit floating point EXR format, and shown tone mapped here for display (with $\gamma = 1.6$).

 a state-of-the-art method for single-day outdoor PS with weakly-calibrated lighting, which is specifically designed to work on the ambiguous case of clear days. The method is robust to shadows, specularities and arbitrary but spatially uniform albedo.

81 Our contributions show that PS can be applied to 82 images obtained over a single day under most weather 83 conditions.

84 2 RELATED WORK

This section focuses on the more relevant work on outdoor PS, for conciseness. For an overview of general PS, the reader is referred to the recent, excellent review in [10].

Woodham's seminal work [1] shows that, for Lambertian 88 surfaces, calibrated PS computes a (scaled) normal vector in 89 90 closed form as a simple linear function of the input image pixels; this linear mapping is only well-defined for images 91 92 obtained under three or more (known) non-coplanar lighting directions. Subsequent work on outdoor PS has strug-93 gled to meet this requirement since, over the course of a 94 day, the sun shines from directions that nearly lie on a 95 plane. These co-planar sun directions then yield an ill-96 posed problem known as two-source PS; despite extensive 97 research using integrability and smoothness constraints [11], 98 [12], results still present strong regularization artifacts on 99 surfaces that are not smooth everywhere. To avoid this issue 100 in outdoor PS, authors initially proposed gathering months 101 of data, watching the sun elevation change over the sea-102 sons [5], [6]. Shen et al. [13] noted that the intra-day copla-103 narity of sun directions actually varies throughout the year, 104 with single-day outdoor PS becoming more ill-posed at 105 high latitudes near the winter solstice, and worldwide near 106 107 the equinoxes.

108 Another important issue is that, so far, most of the literature on PS has adopted a simple directional illumination 109 model, for which optimal lighting configurations can be the-110 oretically derived [13], [14], [15]. Until recently, no attempt 111 had been made to model natural lighting more realistically 112 in an outdoor setup, where lighting cannot be controlled 113 and atmospheric effects are difficult to predict. In such an 114 uncontrolled environment, exploiting the subtlety and rich-115 ness of natural lighting is key to improve the conditioning 116 of PS and successfully apply it in the wild and with short 117 intervals of time. 118

Thus, more recent approaches have attempted to compensate for limited sun motion by adopting richer illumination models that account for additional atmospheric factors in the sky. This is done by employing (hemi-)spherical environment maps [16] that are either real sky images [4], [8], [17], 123 [18] or synthesized by parametric sky models [19], [20]. 124 Despite these developments, state-of-the-art approaches in 125 calibrated [8] and semi-calibrated [20] (based on precise geolocation) outdoor PS are still prone to potentially long waits 127 for ideal conditions to arise in the sky; and verifying the occurrence of such events is still a trial-and-error process. 129

Under more extreme ambiguity, techniques for shapefrom-shading (SfS) [21], [22], [23] attempt to recover 3D normals from a single input image, in which case the shading cue alone is obviously insufficient to uniquely define a solution. Thus, SfS relies strongly on priors of different complexties and deep learning is quickly bringing advances to the field [24], [25], [26]. While this is encouraging, here we seek to improve the accuracy of 3D normal estimation by relying less heavily on priors and more strongly on the photometric cues obtained from multiple images. Finally, most of these methods are limited to a specific type of object and reflectance model (*e.g.*, human faces, Lambertian [25]).

3 OVERVIEW

In this paper, we investigate the complex, natural lighting 143 phenomena that help condition outdoor PS. Our analysis 144 uses the Laval HDR Sky Database [9], [27], a rich dataset of 145 high dynamic range (HDR) images of the sky, captured 146 under a wide variety of weather conditions. In all, the dataset totals more than 5,000 illumination conditions, captured 148 over 50 different days. Fig. 1 shows examples of these environment maps, which are tone mapped for display only; the 150 actual sky images have a dynamic range that spans the full 22 stops required to properly capture outdoor lighting.

142

Our investigations have approached outdoor PS under 153 two different scenarios, leading to two specialized solutions: 154

 First, we consider the popular case of calibrated, outdoor PS for Lambertian objects (Section 4) and we assess how outdoor PS is conditioned solely by the few photometric cues obtained over the course of ne day. By considering Lambertian reflectance, the number of unknowns is reduced to a minimum and, therefore, this analysis provides an upper bound on 161

76

77

78

79

the quality of recovered normals for objects with 162 general reflectances. As we initially reported in [28], 163 [29], partly cloudy days are in fact better for single-164 day outdoor PS since clouds obscure and further 165 scatter sun light, causing a beneficial shift in the 166 effective direction of illumination. Such conditions 167 lend themselves well to calibrated PS algorithms. On 168 the other hand, our analysis also suggest that a dif-169 ferent approach is needed for outdoor PS with clear 170 skies and objects with more general reflectances. 171

2) Second, we consider non-Lambertian objects and 172 the more difficult, under-constrained case of sunny 173 days with clear skies (Section 5). In addition, we 174 also relax the assumption on fully-calibrated light-175 ing. Since there are more unknowns in this new sce-176 177 nario, we cannot rely solely on the few photometric cues obtained within a single day. We thus propose 178 179 an approach that uses deep learning to resolve ambiguities in outdoor PS by aggregating prior 180 181 knowledge on object geometry, material and their interaction with natural outdoor illumination. This 182 new approach is the first of its kind—so far, deep 183 PS had only been applied in indoor scenarios with 184 rich and controlled illumination [10], [30], [31], [32]. 185 We conclude by discussing how the advantages of the 186

two solutions above could be integrated into a single, more
 generic approach in future work.

189 4 LAMBERTIAN, CALIBRATED OUTDOOR PS

190 4.1 Image Formation Model

Consider a small, Lambertian surface patch with normal vector **n** and albedo ρ (w.l.o.g., assume albedo is monochromatic). At time *t*, this surface patch is observed under natural, outdoor illumination represented by the environment map $L_t(\omega)$ (e.g, Fig. 1), with ω denoting a direction in the unit sphere. With an orthographic camera model, this patch is depicted as an image pixel with intensity

$$b_t = \frac{\rho}{\pi} \int_{\Omega_{\mathbf{n}}} \mathbf{L}_t(\boldsymbol{\omega}) \langle \boldsymbol{\omega}, \mathbf{n} \rangle d\boldsymbol{\omega} , \qquad (1)$$

where $\langle \cdot, \cdot \rangle$ denotes the dot product. Integration is carried 200 out over the hemisphere of incoming light, Ω_n , defined by 201 the local orientation n of the surface, Fig. 2. This hemisphere 202 corresponds to an occlusion (or attached shadow) mask; 203 204 only half of the pixels in the environment map contribute to the illumination of the surface patch. To make the analysis 205206 tractable and independent of object geometry, this analysis focuses on the simpler case without cast shadows. 207

This image formation model is then discretized as,

$$b_t = \frac{\rho}{\pi} \sum_{\boldsymbol{\omega}_j \in \Omega_{\mathbf{n}}} \hat{\mathbf{L}}_t(\boldsymbol{\omega}_j) \langle \boldsymbol{\omega}_j, \mathbf{n} \rangle , \qquad (2)$$

210

216

208

199

with $\mathbf{L}_t(\boldsymbol{\omega}_j) = \mathbf{L}_t(\boldsymbol{\omega}_j) \Delta \boldsymbol{\omega}_j$ representing the environment map weighted by the solid angle $\Delta \boldsymbol{\omega}_j$ spanned by pixel j ($\Delta \boldsymbol{\omega}_{jr}, \forall j$, are normalized as to sum to 2π). Eq. (2) can be further summarized into the equivalent form

$$b_t = \overline{\mathbf{I}}_t^T \mathbf{x},\tag{3}$$



217



Fig. 2. A normal **n** defines an integration hemisphere Ω_n on the environment map. Only light emanating from this hemisphere contributes to the shading on that patch. Thus, patches with different normals are lit differently even if the environment map is the same.

where $\mathbf{x} = \rho \mathbf{n}$ is the albedo-scaled normal vector and

$$\bar{\mathbf{l}}_t = \frac{1}{\pi} \sum_{\boldsymbol{\omega}_j \in \Omega_n} \hat{\mathbf{L}}_t(\boldsymbol{\omega}_j) \boldsymbol{\omega}_j \in \mathbb{R}^3.$$
(4)

This vector $\overline{\mathbf{l}}_t$ can be interpreted as a virtual point light 220 source summarizing the illumination provided by the envi-221 ronment map at a time t. This vector $\overline{\mathbf{l}}_t$ is the mean of the 222 light vectors computed over the hemisphere of incoming 223 light directions defined by \mathbf{n} (see Fig. 6). As such, this vector 224 is henceforth referred to as the *mean light vector* (MLV). It is 225 important to note that, as opposed to the traditional PS scenario where point light sources are fixed and thus indepen-227 dent of \mathbf{n} , here the *per-pixel* MLV is a function of \mathbf{n} . Thus, 228 patches with different orientations define different sets of 229 MLVs (see Fig. 2). A similar lighting representation has 230 been adopted in the uncalibrated case in [18].

Given multiple images taken at times $t \in \{1, 2, ..., T\}$, 232 we collect all photometric constraints for patch x to obtain: 233

$$\mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_T \end{bmatrix} = \begin{bmatrix} \mathbf{\overline{I}}_1^T \\ \mathbf{\overline{I}}_2^T \\ \vdots \\ \mathbf{\overline{I}}_T^T \end{bmatrix} \mathbf{x} = \mathbf{L}\mathbf{x} \,. \tag{5}$$

With Eq. (5), this model of natural environmental illumina- 236 tion becomes quite similar to a model with a distant point 237 light source, the well-known case in PS. However, note that 238 each \bar{l}_t in L is a function of Ω_n and, thus, of n. 239

Most importantly, in outdoor PS, a well-defined solution 240 x may exist even if the relative sun motion is nearly planar 241 during a certain time interval. Instead of relying solely on 242 sun direction, now, the solution requires non-coplanar 243 mean light vectors \overline{l}_t , which are determined by a compre-244 hensive set of natural illumination factors. 245

4.2 Measuring Uncertainty

From Eq. (5), the least-squares solution $\mathbf{x} = (\mathbf{L}^T \mathbf{L})^{-1} \mathbf{L}^T \mathbf{b}$ of 247 outdoor PS is clearly affected by the condition number of \mathbf{L} . 248 Thus, we next characterize how well the solution \mathbf{x} is con-249 strained by natural, outdoor illumination within a given 250 time interval (e.g, one day)—which is encoded by the set of 251 mean light vectors $\mathbf{\bar{l}}_t$ in \mathbf{L} or, equivalently, the set of environ-252 ment maps $L_t(\cdot)$. 253

To assess the reliability of a solution x, we follow stan- 254 dard practice in PS [15], [33] and consider image measure- 255 ments corrupted by zero-mean Gaussian noise with equal 256

261

267

273

274

281

296

300

variance σ^2 (as least squares estimation is only optimal for this practical, most common noise model). Thus, b in Eq. (5) follows a normal distribution:

$$\mathbf{b} \sim \mathcal{N}(\boldsymbol{\mu}_b, \sigma^2 \mathbf{I}) \,, \tag{6}$$

where μ_b has the (unknown) uncorrupted pixel values.

Since the desired least-squares solution for the albedoscaled normal, $\mathbf{x} = (\mathbf{L}^T \mathbf{L})^{-1} \mathbf{L}^T \mathbf{b}$, is a linear transformation of a Gaussian random vector, it is easy to show that

 $\mathbf{x} \sim \mathcal{N}\left(\boldsymbol{\mu}_x, \sigma^2 (\mathbf{L}^T \mathbf{L})^{-1}\right),$ (7)

where $\mu_x = (\mathbf{L}^T \mathbf{L})^{-1} \mathbf{L}^T \mu_b$ is the expected value of x. Once the albedo of a surface patch is known, we analyze its contribution to the uncertainty in the estimated normal vector, $\mathbf{n} = \rho^{-1} \mathbf{x}$, using a similar distribution,

$$\mathbf{n} \sim \mathcal{N}\left(\frac{\boldsymbol{\mu}_x}{\rho}, \frac{\sigma^2}{\rho^2} \left(\mathbf{L}^T \mathbf{L}\right)^{-1}\right).$$
 (8)

The marginal distributions in Eq. (8) allow us to derive confidence intervals that indicate the uncertainty in each component of the least squares estimate $\hat{\mathbf{n}} = [\hat{n}_x \ \hat{n}_y \ \hat{n}_z]^T$ of $\mathbf{n} = [n_x \ n_y \ n_z]^T$. The corresponding 95 percent confidence interval [34] is given by

$$\hat{\mathbf{n}} \pm \boldsymbol{\delta}, \quad \text{with } \boldsymbol{\delta}_k = 1.96 \frac{\sigma \lambda_k}{\rho}, \qquad (9)$$

where λ_k is the square root of the *k*th element on the diagonal 282 of $(\mathbf{L}^T \mathbf{L})^{-1}$. As expected, the sensor-dependent noise level σ 283 is not the only factor that determines uncertainty. The noise 284 gain factor λ_k in Eq. (9) reveals how outdoor illumination 285 (the conditioning of L) can amplify the effect of noise on the 286 287 solution $\hat{\mathbf{n}}$. The albedo ρ also impacts the solution stability, where a lower albedo translates in a larger variance in the 288 obtained estimate $\hat{\mathbf{n}}$ (as less light is reflected towards the 289 camera). Our goal is then to answer the remaining question: 290 how do natural changes in outdoor illumination affect this 291 noise gain factor (λ_k) and, therefore, uncertainty? 292

To provide a measure of uncertainty that is more intuitive than Eq. (9), we consider angular distances in degrees,

$$\theta^{\pm} = \cos^{-1}(\mathbf{n}^T \hat{\mathbf{n}}^{\pm}), \qquad \hat{\mathbf{n}}^{\pm} = \frac{\hat{\mathbf{n}} \pm \delta}{\|\hat{\mathbf{n}} \pm \delta\|}.$$
(10)

The uncertainty in the estimate of n is then summarized in a single confidence interval, in degrees,

$$\mathcal{C}_{\mathbf{n}} = \left[0, \, \max(\theta^{\pm}) \, \right], \tag{11}$$

which indicates the expected accuracy of the estimated surface orientation $\hat{\mathbf{n}}$.

Note that the condition number, determinant, and trace 303 of matrix $(\mathbf{L}^T \mathbf{L})^{-1}$ can also be used as measures of total vari-304 ance in the estimated solutions-as done in [33]-to find the 305 optimal location of point light sources in PS. These meas-306 ures are closely related to the rank of matrix L, which must 307 be three for a solution to exist; that is, $L^{T}L$ must be nonsin-308 gular. In practice, this matrix is always full-rank, although it 309 is often poorly conditioned [13]. In the following, we also 310 consider the gain factor λ_k in (9) as a measure of uncertainty 311



Fig. 3. Median confidence interval of normal estimates (red line) as a function of mean sun visibility over the course of the day for a signal noise $\sigma = 0.5\%$, in bins of 10° . Our analysis predicts that normal reconstruction errors will likely be high if the sky is completely overcast (low sun visibility), or completely clear (high sun visibility). Good results can thus be expected in partially cloudy conditions, as shown in Fig. 4. The lower (upper) edge of each blue box indicates the 25th (75th) percentile. Statistics are computed only on normals pointing upwards to lessen ground effects.

independent of albedo and sensor noise. We focus on analyzing our ability to recover geometry and will assume that the albedo is constant. 314

315

4.3 Effect of Clouds on Outdoor Lighting

This section investigates the environmental element that ³¹⁶ most influences mean light vectors throughout the day: ³¹⁷ *clouds*. Cloud coverage has an important effect on the uncertainty of normal reconstruction because clouds introduce ³¹⁹ variability in illumination and, thus, new photometric cues ³²⁰ as they (dis)occlude the sun. Here, we present a systematic ³²¹ analysis of their influence on outdoor PS. To control for the ³²² effect of the sun elevation, the analysis is performed on ²³ ³²³ days with similar sun elevations by keeping only the skies ³²⁴ captured in October and November. ³²⁵

We approximate cloud coverage by computing the fraction of time that the sun is visible, i.e, that it fully shines on the scene, for a given day. To do so, we simply find the series spot in a sky image, and determine that the sun shines on the scene if the intensity of the brightest pixel is greater than 20 percent of the maximum sun intensity—we determined empirically that this is the point at which the sun is bright enough to start creating cast shadows. Cloud coverage is represented by computing the mean sun visibility for a given day. Values of less than 10–15 percent indicate mostly overcast skies, while skies are mostly clear if above 85–90 percent.

The relation between sun visibility and the confidence ³³⁸ interval C_n is shown in Fig. 3. Photometric-related normal ³³⁹ reconstruction errors will likely be quite high in two situa- ³⁴⁰ tions: completely overcast (low sun visibility), or completely ³⁴¹ clear skies (high sun visibility). Interestingly, good recon- ³⁴² struction results are expected for a wide range of cloud cov- ³⁴³ erage conditions, ranging from 10–90 percent mean sun ³⁴⁴ visibility. ³⁴⁵

These results are corroborated by Fig. 4, which shows the 346 confidence intervals themselves. These intervals are aver- 347 ages over skies belonging to four groups: overcast (0–15 per- 348 cent), mixed overcast (15–50 percent), mixed clear (50–85 349 percent), and clear (85–100 percent) days. Again, high 350



Fig. 4. Influence of cloud cover on the 95 percent confidence intervals (in degrees) with $\sigma = 1\%$. Each pair of plots show the full sphere of normals from two different viewpoints: South (left), and North (right). Four different types of skies are shown, based on sun visibility. For example, the top-left plots show the confidence intervals averaged over all days with direct sun visibility in the range 85-100 percent.

uncertainty results are visible for the two extreme cases of
fully overcast and fully clear days, while the remainder
indicate more stable solutions.

Under clear skies, the MLVs \overline{I}_t of the model above will 354 point nearly towards the sun, from which arrives most of 355 the incoming light. Thus, near an equinox (worldwide), the 356 associated MLVs are nearly coplanar [13], resulting in poor 357 performance, Fig. 5a. For a day with an overcast sky, perfor-358 359 mance is also poor because the set of MLVs are nearly colinear and shifted towards the patch normal n, Fig. 5b. The 360 improved conditioning in mixed skies is explained by the 361 following key observation: cloud cover *shifts* the MLVs \bar{l}_t 362 towards zenith and away from sun trajectory in the sky, 363 Fig. 5c. Therefore, even when the sun moves along a trajec-364 tory that nearly lies on a 3D plane, as shown in Fig. 6, cloud 365 cover effectively causes an out-of-plane shift of the MLVs, 366 making reconstruction possible. 367

It is important to note that surface patches of different 368 orientations (normals) are exposed to different hemispheres 369 of illumination, with light arriving from above (sky) and 370 below (ground). More MLV trajectories are shown in Fig. 7 371 372 for three different normal vectors (rows) and two different days (columns). Each globe represents the coordinate sys-373 374 tem for the environment maps captured in a day. For each combination normal-day, the time-trajectory of computed 375 376 MLV directions (dots) and intensities (colors) are shown on the globe. Brighter MLVs lie close to the solar arc, while 377 darker MLVs may shift away from it. Note that we present 378 normals that are mainly Southward as they receive the most 379 direct sunlight throughout the day in the Northern hemi-380 sphere. Surfaces with normals pointing North, for example, 381 would be in shadow throughout the day in latitudes higher 382



Fig. 5. Impact of cloud coverage on the numerical conditioning of outdoor PS: Clear (a) and overcast (b) days present MLVs with stronger coplanarity; in partly cloudy days (c) the sun is often obscured by clouds, which may lead to out-of-plane shifts of MLVs.



Fig. 6. Cloud effect on the MLV over one day: while the sun path (orange) yields nearly co-planar directions of illumination, the mean light directions (red dots) for a normal pointing up provide a much more varied set (data from 11/06/2013, second row of Fig. 1).

than the Tropic of Cancer around the winter solstice. Thus, 383 the remainder of this paper considers a camera pointing 384 North. 385

To more closely match the scenario considered above, we 386 scale these real MLVs so that the brightest one over all days 387 (i.e, for the most clear sky) has unit-length. From Fig. 7, also 388 note that some MLVs are shifted very far from the solar arc 389 but, as indicated by the darker colors, their intensity is 390



(a) Mixed clouds (06-NOV-13) (b) Mixed clouds (11-OCT-14)

Fig. 7. Globes representing the coordinate system of sky probes. Each normal (blue arrow) defines a shaded hemisphere in the environmental map that does not contribute light to the computed MLVs (dots). All MLVs in two particular partly cloudy days (columns) were computed from real environment maps [28] for 3 example normal vectors (rows). Relative MLV intensities are shown in the color bar on the left.



Fig. 8. Noise gain for normal directions n of patches visible to the camera, which is located South of the hypothetical target object. The colors indicate the shifting (coplanarity) of the associated MLVs. On both days, normals that are nearly horizontal are associated with nearly coplanar MLVs (smaller shifts, higher gains). These normals define a *zerocrossing region* between positive and negative out-of-plane shifts (mid row in Fig. 7), where sun occlusion shifts MLVs predominantly *along* the solar arc.

dimmed considerably by cloud coverage; little improve-ment in conditioning is obtained from these MLVs.

Most important, Fig. 7 shows that the amount of out-of-393 plane MLV shift (elevation) relative to the solar arc also 394 depends on the orientation n of the surface patch. This indi-395 cates that outdoor PS may present different degrees of uncer-396 tainty (conditioning) depending on the normal of each patch. 397 Indeed, the maximum noise gain ($\lambda_{max} = max(\lambda_k)$) values in 398 Fig. 8 show that patches with nearly horizontal normals 399 (orthogonal to the zenith direction) are associated with sets 400 of MLVs that are closer to being coplanar throughout the 401 day. As expected, patches oriented towards the bottom also 402 present worse conditioning since they receive less light. 403

This key observation also demonstrates the advantages 404 of adopting more elaborate illumination models (e.g, [8]). 405 For instance, the simpler point light model was used in [13] 406 to study the conditioning of outdoor PS. Because the atmo-407 spheric component is not modeled, the conclusion was that 408 single-day reconstruction breaks down in two cases of 409 nearly coplanar sun directions: closer to the poles near the 410 winter solstice, and worldwide near an equinox. Our results 411 suggest that more attention should be placed on the illumi-412 nation model, without focusing exclusively on the sun. 413

414 4.4 Lambertian, Calibrated PS on Partially Cloudy Days

The analysis performed on the HDR sky dataset (c.f. Section 3)
indicates that surface patches may be better reconstructed in
certain conditions, dependent upon cloud coverage and the

orientation of the patch itself. In the case of partially cloudy 419 days, our investigation reveals that those conditions usually 420 shift the MLVs enough for outdoor PS methods to work. 421

To validate that accurate surface reconstructions can 422 indeed be obtained on partially cloudy days, we captured a 423 sequence of a real object lit by the sky over the course of a 424 day. We oriented an owl statuette towards south and took 425 66 HDR captures using a Canon EOS Rebel SL1 between 426 10h30 and 16h30, local time, in Quebec City. We simulta-427 neously captured hemispherical HDR sky images (as in Section 3) to provide high fidelity estimates of the illumination 429 conditions for each image as shown in Fig. 9. Ground-truth 430 surface normals were obtained by aligning a 3D model of 431 the object (obtained with a Creaform MetraSCAN scanner) 432 to the image using POSIT [35].

We then perform calibrated outdoor PS on these images 434 using the algorithm proposed by Yu *et al.* [8], with the fol- 435 lowing three differences: (*i*) we use all possible pairs of 436 images to compute ratios, instead of selecting a single 437 denominator image; (*ii*) we apply anisotropic regulariza- 438 tion [12] to mitigate the impact of badly-conditioned pixels 439 on the surface reconstruction; and (*iii*) remove the low-rank 440 matrix completion preprocessing, which, in our experi- 441 ments, caused slightly degraded performance. 442

The result on these real images is shown on Fig. 9. As 443 predicted in the analysis from Fig. 8, normals on the head 444 and the bottom of the abdomen, pointing respectively up or 445 down, are mostly accurately estimated. As can be observed 446 in Fig. 9a, clouds sometimes occlude the sun, which 447 improves the conditioning of the problem to yield an 448 acceptable result. Without clouds, this day would have lead 449 to an unstable formulation of the photometric stereo prob-450 lem, as it is close to the fall equinox, which corresponds to 451 the worst case scenario with coplanar sun directions [13].

5 NON-LAMBERTIAN PS ON CLOUDLESS DAYS

453

Full environment maps taken at short intervals are needed 454 to analyze the case of partly cloudy days, as one needs to 455 capture the precise moment when the sun is occluded by 456 clouds. However, in the case of clear days, the photometric 457 cue is much weaker but the general appearance of the sky is 458 more predictable and can be modeled by physically-based 459 sky models (as in [20]). In this section, we present a novel 460 approach using deep learning to handle the ambiguities 461 that arise in outdoor PS on a single day with a clear sky.



Fig. 9. (left) Real outdoor HDR images of owl statuette and corresponding HDR environment maps (top row) providing synchronized, high-fidelity estimates of illumination conditions. All images were acquired on 10/11/2014 and tone-mapped for display only (with $\gamma = 1.6$). The sun visibility was 43 percent on this day. We show the ground truth normals of the object (a) as well as normals recovered from [8] (b), along with a reference normal sphere in inset. The reconstruction error (c) shows sphere is shown as a color coding reference; (b) normal estimation error at each pixel; and (c) the error distribution, in degrees.



Fig. 10. Our novel CNN architecture for deep single-day outdoor PS on sunny days. The network operates on 16×16 patches B_t of the input image, captured at 8 time intervals *t* regularly spaced throughout a single day. The network uses convolutional (blue) and residual (red) layers before estimating the normals using fully-connected layers (green). Two losses are used to train our method, one based on the cosine distance with the ground truth \mathbf{f} and another to constrain the norm of the output vector.

Our CNN-based approach compensates for the lack of 463 photometric constraints by modeling prior knowledge on 464 object geometry, material properties, as well as their local 465 spatial correlation and interaction with natural outdoor 466 lighting. In order to build such knowledge base, one needs a 467 large number of images depicting various objects lit by out-468 door lighting throughout the day, over different geographic 469 locations and days over the year; finally, the surface normal 470 map of each object is also required. Unfortunately, no such 471 large-scale dataset currently exists, so a natural choice is to 472 synthesize realistic data to train our network. Next, we pres-473 ent our problem formulation, CNN architecture, followed by 474 475 the training procedure and data generation.

476 5.1 Illumination Model: The Solar Analemma

We follow a semi-calibrated PS approach that does not 477 require known lighting environments [8] nor complete cam-478 479 era geolocation data [20]. Our method only assumes that: (1) the object images are captured at approximate prede-480 fined times of the day, $t \in \{t_1, t_2, \dots, t_T\}$; (2) the sun is unob-481 structed by clouds at these times; and (3) the camera is 482 orthographic and faces approximately North (or South). In 483 Section 5.6, we analyze the robustness of our network with 484 respect to departures from these ideal conditions. 485

Together, these assumptions constrain the sun position to 486 lie within an "8-figure" subspace at each time t, known as a 487 solar analemma, whose shape also varies with geographical 488 location (Fig. 11). For a given time t, the sun may be positioned 489 490 at different locations depending upon the selected date and latitude, as prescribed by the analemma. The neural network 491 is thus expected to adapt to this (constrained) variability in 492 sun position and associated intensity. As shown in Figs. 11a 493 and 11b, for a given timestamp and latitude, the sun position 494 495 spans relatively small angular ranges, which still remain quite constrained even when considering geographical locations 496 497 sampled over the Northern Hemisphere (Fig. 11c) (note that a similar plot would be obtained by sampling the Southern 498 Hemisphere with the camera facing South). 499

Clear days can be accurately synthesized by parametric sky models, with much lower dimensionality in comparison to a full environment map. To generate training data, we use the physically-based parametric sky model of Hošek and Wilkie [36] to obtain the spherical illumination function $L_t(\omega)$ in Eq. (1). The model represents the spectral sky radiance as a parametric function of the sun position, sky turbidity and 506 ground albedo; turbidity is set to 2, which corresponds to a 507 clear day, and ground albedo to 0.3. Note that we do not 508 model light scattering caused by clouds obscuring the sun 509 and thus assume the sun is fully visible in the sky. 510

5.2 Deep Outdoor PS Network

Here, we consider a more general image formation model in 512 which the Lambertian term $\frac{\rho}{\pi}$ in Eq. (1) is replaced with a 513 standard GGX shader $\rho(\cdot)$ with varying diffuse and specular 514 parameters. Our goal now is to invert this new rendering 515 equation and recover the surface normal **n** based on the 516 observed changes in pixel intensities b_t , which are caused 517 by the changing natural illumination $\mathbf{L}_t(\boldsymbol{\omega})$ throughout the 518 day. However, a solution based solely on the photometric 519 cues from a sunny day is typically undefined due to limited 520 sun motion and, thus, insufficient variability in $\mathbf{L}_t(\boldsymbol{\omega})$ and b_t . 521

Therefore, instead of considering a single pixel b_t , we 522 reformulate our goal and instead aggregate additional RGB 523 image data within a *neighborhood* $\mathbf{B}_t \in \mathbb{R}^{P \times P \times 3}$, depicting a 524 larger surface patch of width P centered at the pixel \mathbf{b}_t . 525 Now we seek to learn a predictor $\mathbf{N} = f(\mathbf{B}_{t_1}, \dots, \mathbf{B}_{t_T})$, 526 where T denotes the number of input images and 527 $\mathbf{N} \in \mathbb{R}^{P \times P \times 3}$ is the patch normals. In this paper, T = 8 and 528 P = 16 but we experiment with other values in Sec- 529 tions. 5.6.2 and 5.6.3 respectively. This approach is motivated by the fact that complex object geometry is often 531 made up of simpler, small surface patches presenting highly 532 correlated surface normals and material properties. 533

A natural way to obtain this predictor $f(\cdot)$ is to train a 534 Convolutional Neural Network (CNN) and learn a nonlin-535 ear function of local surface features that are highly correlated with the normal n at the center of the patch. We train 537 our network on a large synthetic database of surface patches realistically rendered (Section 5.3). 539

5.2.1 Network Architecture

The function $\mathbf{N} = f(\mathbf{B}_{t_1}, \dots, \mathbf{B}_{t_T})$ introduced above is 541 designed as CNN with the architecture shown in Fig. 10. The 542 network takes 8 input 16 × 16 image patches, extracted from 543 8 images captured at regular intervals Δt between 9:00 and 544 16:00 solar time throughout a single sunny day. Note that no 545 other information (geolocation, capture time, etc.) is pro-546 vided to the network. The first layer is composed of 32 547

511



Fig. 11. Solar analemma: position of the sun in the sky at a specific time of the day and throughout a year over (a) Paris and (b) the Tropic of Cancer. Note how the analemmas spread over a wide range of zenith and azimuth angles over the course of a year. (c) Probability of the sun location in the sky for our training set.

channels of 5×5 filters with shared weights across the 8 548 inputs. The resulting feature maps are subsequently 549 concatenated in a single $14 \times 14 \times 256$ feature tensor. A sec-550ond convolutional layer is then used, yielding 256 channels, 551 followed by 3 residual blocks as defined in the resnet-18 552 architecture [37]. Lastly, 2 fully-connected layers (FC) are 553 used to produce a $16 \times 16 \times 3$ patch of estimated normals **n**. 554 Note that we experimented with fully-convolutional archi-555 tectures [32] but found the FC layers to yield better results. 556 The ELU activation function [38] is used at every convolu-557 tional and fully connected layer, except the output layer 558 where a $tanh(\cdot)$ function is used. As in [39], batch normaliza-559 tion [40] is applied at every layer except the first and the out-560 put laver. 561

The 16×16 estimated normals are represented by Cartesian (x, y, z) components of the surface normal of the input patch. We experimented with parameterization in both Cartesian (x, y, z) and spherical coordinates (θ, ϕ) , but found the former to be more stable despite its additional degree of freedom. We hypothesize this may be due to the "wraparound" issue with the azimuth angle ϕ .

To process entire images, we crop overlapping tiles from the image with a stride of 8 pixels. Since a pixel can belong to up to 4 patches, the network produces several estimates **f** that are then merged together using a weighted average. We use a Gaussian kernel with $\sigma = 4$ centered on the middle of the patch as weighting function to perform the linear interpolation across overlapping patches.

576 5.2.2 Training

585

589

The network learns a function that estimates the patch normals **N**. We define the loss to be minimized between the estimated and ground truth patch normals **N** and **N**^{*} respectively as the sum of two separate loss functions defined on individual patch normals \mathbf{n}_i , $i \in \{1, ..., N\}$ where $N = 16 \times 16 = 256$. The total loss is the sum over all *N* individual normals:

$$\mathcal{L}(\mathbf{N}, \mathbf{N}^*) = \sum_{i=1}^{N} \left(\mathcal{L}_{\cos}(\mathbf{n}_i, \mathbf{n}_i^*) + \mathcal{L}_{\text{unit}}(\mathbf{n}_i) \right).$$
(12)

The first term is the cosine distance between the estimated \mathbf{n}_i and ground truth normal \mathbf{n}_i^* :

$$\mathcal{L}_{\cos}(\mathbf{n}_i, \mathbf{n}_i^*) = 1 - \frac{\left\langle \mathbf{n}_i, \mathbf{n}_i^* \right\rangle}{\|\mathbf{n}_i\| \|\mathbf{n}_i^*\|}, \qquad (13)$$

where $\langle \cdot, \cdot \rangle$ denotes the dot product. The second term enfor- 590 ces the unit-length constraint on the recovered normal: 591

$$\mathcal{L}_{\text{unit}}(\mathbf{n}_i) = |\|\mathbf{n}_i\| - 1|.$$
 (14)

601

The loss in Eq. (12) is minimized via stochastic gradient ⁵⁹⁴ descent using the Adam optimizer [41] with an initial learn- ⁵⁹⁵ ing rage of $\eta = 0.001$, a weight decay $\lambda = 1 \times 10^{-4}$ and the ⁵⁹⁶ recommended values $\beta_1 = 0.9$ and $\beta_2 = 0.999$. Mini-batches ⁵⁹⁷ of 128 samples were used during training and regularized ⁵⁹⁸ via early stopping. Training typically converges in around ⁵⁹⁹ 250 epochs on our dataset, which is described next.

5.3 Training Dataset

To train our predictor function $f(\cdot)$, we rely on a large training 602 dataset of synthetic objects, lit by a physically-based outdoor 603 daylight model. To generate a single 8-images set of inputs, 604 we randomly select a combination of: 1) object shape, 2) material, and 3) geo-temporal coordinates for lighting. We now 606 detail how each of these 3 choices are made. 607

Since the neural network only sees patches of 16×16 608 pixels, its receptive field is, by design, not large enough to 609 learn priors on whole object shapes. Therefore, our dataset 610 contains a wide variety of local surface patches. We used 611 the blob dataset from [22] as training models. We also 612 added simple primitives (cube, sphere, icosahedron, cone) 613 to the data. A validation set, comprised of one of the blobs 614 models that was kept from the training set as well as some 615 models from the Stanford 3D Scanning Repository [42] and 616 the owl model used in [28], was also created. All blobs and 617 geometric primitives are randomly rotated about their 618 centroid.

To model a wide range of surface appearance ranging 620 from diffuse to glossy, we employ a linear combination of a 621 lambertian and a microfacet model: 622

$$\rho(\boldsymbol{\omega}, \mathbf{v}, \mathbf{n}) = \boldsymbol{\rho}_c(\boldsymbol{\alpha} + (1 - \boldsymbol{\alpha})\rho_{\text{GGX}}(\boldsymbol{\omega}, \mathbf{v}, \mathbf{n}, \sigma)), \qquad (15)$$

where $\rho_c \in \mathbb{R}^3$ is the surface color, and ρ_{GGX} is the GGX 625 microfacet model [43] with surface roughness parameter σ . 626

The albedo ρ is generated in HSV space, where 627 $H \sim U(0,1)$, $S \sim T(0,0,1)$, and $V \sim T(0,0.75,1)$, where 628 U(a,b) is a uniform distribution in the [a,b] interval and 629 T(a,b,c) is a triangular distribution in the [a,c] interval with 630 mode b. This generates colors that are in general bright and 631 prevents an abundance of strongly saturated colors. Surface 632 roughness σ is sampled as $\sigma \sim T(0.2, 0.4, 1)$ to avoid mirror- 633 like surfaces. Finally, we sample a mixing coefficient 634 $\alpha \sim U(0, 1)$.

To light the scene with a wide variety of realistic outdoor 636 lighting conditions, we rely on the Hošek-Wilkie physically- 637 based sky model [36]. We also placed a ground plane of 638 albedo 0.3 outside the field of view of the camera, to gener- 639 ate a light bounce from below the object. 11 random loca- 640 tions in the Northern Hemisphere between latitude 0° 641 (Equator) and $^{56°}$ (Moscow) were selected. Furthermore, 6 642 random days throughout the year were chosen in addition 643 to the equinoxes and solstices. This results in 110 pairs of 644 geographical locations and dates, which are used to com- 645 pute the sun position in the sky throughout the day using 646 [44]. The distribution of the resulting sun positions through 647 out our training set is shown in Fig. 11. For every pair of 648



Fig. 12. (top) An example of the lighting environment maps and renders throughout a day. (bottom) Qualitative results (odd rows) and errors in degrees (even rows) of our technique and the state-of-the-art on single-day photometric stereo in the semi-calibrated [20] and calibrated [8] cases, deep photometric stereo [31] and single image normal estimation [24], [26] (averaged over the day) on our real lighting dataset. *More results available in the supplementary material, which can be found on the Computer Society Digital Library at http://doi.ieeecomputersociety. org/TPAMI.2019.2962693.*

geographical location and day, 8 timestamps ranging from 649 9:00 to 16:00 are used to perform the renders. Timestamps 650 are aligned to the solar noon instead of the political time 651 652 zone of the geographic location. Although we sample only geographical locations in the Northern hemisphere, our 653 654 dataset represents equally well days in the Southern hemisphere. Indeed, flipping the images left-right, reversing the 655 image order (from 16:00 to 9:00) and pointing the camera 656 Southward would generate data identical to our training 657 dataset. 658

The resulting images are rendered with the Cycles physi-659 cally-based rendering engine. This results in a dataset of 660 369,440 renders corresponding to 23,090 combinations of 661 geo-temporal coordinates and materials properties, which 662 we then split into 21,220 and 1870 for training and valida-663 tion, respectively. Each render has a resolution of 256×256 664 pixels, which amounts to over 10 millions input-output 665 pairs of 16×16 patches to train on. Special care was taken 666 into ensuring no 3D model nor material properties were 667 shared between both the training and validation datasets. 668 Please see the supplementary material for example training 669 images, available online. 670

671 5.4 Results on Synthetic Images With Real Lighting

We first evaluate and compare the techniques using a dataset of synthetic objects lit by real skies. To generate



Fig. 13. Median reconstruction error on our real lighting dataset displayed vertically as box-percentile plots [46]; the center horizontal bars indicate the median, while the bottom (top) bars are the 25th (75th) percentiles. Our proposed method (green) provides state-of-the-art performance compared to non-learned methods for single-day PS (blue [20], orange [8]), deep learning methods on calibrated photometric stereo (red [31]) and single image normals reconstruction (purple [26], brown [24]).

the images, we manually selected 3 sunny days over 2 674 geographical locations from the Laval HDR sky data- 675 base [9], which contains unsaturated HDR, omnidirec- 676 tional photographs of the sky captured with the 677 approach proposed in [45]. We build a virtual 3D scene 678 containing the HDR sky environment map as the sole 679 light source, a 3D object viewed by an orthographic cam- 680 era, and a 0.3 albedo ground plane placed under the 681 object, outside the field of view of the camera. We used 682 the 3D models from the validation set which the neural 683 network never saw during training. This results in a 684 dataset of 960 renders yielding 60 normal maps to evalu- 685 ate. Example images obtained with this technique are 686 shown in Fig. 12.

We compare our method to several state-of-the-art tech- 688 niques relying on photometric stereo and/or deep learning 689 to estimate surface normals from images. We first compare 690 to the calibrated PS technique from Section 4.4. While it is 691 an improvement over the method of Yu et al. [8], we still 692 refer to it as [8] in figs. 12 and 13. We also compare to the 693 semi-calibrated method of Jung et al. [20], which requires 694 only knowledge of the camera geolocation. For deep learn- 695 ing techniques, we compare to the recent Deep Photometric 696 Stereo Network (DPSN) [31], which operates on one pixel at 697 a time. Since it assumes known point light source lighting, 698 we re-trained this model using the sun position from a geo- 699 graphical location and date representative of our training 700 dataset. In addition, we also compare to single image net- 701 works: Eigen and Fergus [24] and MarrNet [26]. Since they 702 rely on a single image, we take the mean of their results 703 averaged over all 8 inputs. 704

The comparative results, shown qualitatively in Fig. 12 705 and quantitatively in Fig. 13, clearly demonstrate that our 706 approach significantly outperforms all other techniques. We 707 observe that both single image techniques do not work well 708 and result in very high median errors of around 40° and 70° 709 for [26] and [24], respectively. For [24], this is probably due 710 to the fact that they cannot handle the harsh shadows cre- 711 ated by outdoor lighting during sunny days, since they train 712 with indoor lighting only. In addition, MarrNet [26] outputs 713 a voxel occupation grid and only produces normals as a 714



Fig. 14. Result on real statuettes (ill-posed, single day PS problem): (a) example input images around solar noon; (b) the ground-truth (3D-scanned) normals; (c) normals estimated by our method; and (d) angular error map, median error in degree and amount of estimated normals within 30° of the ground truth (R30). The top two rows were taken on 2015-08-22 while the bottom three rows were taken on 2018-05-24.

byproduct (in its latent stage). As such, this method may notbe fully optimized for normal estimation.

The PS techniques yield much better results but still yield 717 quite significant error since sunny days do not contain suffi-718 cient constraints to accurately recover surface normals. The 719 (improved) calibrated method of Yu et al. [8] is comparable 720 to the results obtained by DPSN, with a median normal 721 angular estimation error of 33°. Interestingly, the method of 722 Jung et al. [20] actually yields better results with a median 723 error of 22°, despite needing less information (geolocation 724 and time) than the calibrated methods. This could be due to 725 its reliance on a parametric clear sky model to estimate 726 lighting, which closely matches the actual ground truth 727 lighting, and to its reliance on an intensity profile matching 728 algorithm. 729

Note that most PS techniques capture with some degree 730 731 of success the left/right component of the surface normals (roughly speaking, the red and blue tints in the normal 732 maps). This axis is the same as the sun trajectory through 733 the day when the camera is facing North or South. This 734 results in strong photometric constraints on this axis. On 735 the other hand, the recovery of the up/down axis is much 736 less successful on most techniques as outdoor photometric 737 cues lack information in this direction through a single 738 sunny day. 739

740 In contrast, our method yields a normal map that is, 741 although a bit smoother, qualitatively very similar to the ground truth. Quantitatively, our approach achieves a 742 median error of 14° over the evaluation set, with error predominantly below that of the second best performing 744 method [20] (see Fig. 13). Even when trained on purely synthetic data, our network is able to generalize well to images 746 rendered with real lighting. The difference in performance 747 with respect to DPSN shows the usefulness of dealing with 748 image patches, which allows the network to learn appropriate patch-based shape priors which can be exploited when 750 the photometric cue alone is not sufficient. 751

5.5 Results on Real Captures

We further evaluate our method on real data. We captured 753 sequences of 8 outdoor images of 4 real statuettes during a 754 single sunny day using a tripod-mounted Canon EOS 5D 755 Mark III camera with a 300 mm lens. These HDR images 756 were obtained by merging camera exposure range from 1/757 8000 to 1 second at f/45. Ground-truth normals were 758 obtained from a Creaform GO!SCAN 3D laser scans of the 759 real objects. The results shown in Fig. 14 demonstrate the 760 performance of our proposed method on such ill-posed out- 761 door PS problem. Using photometric cues alone, the two 762 top statuettes from Fig. 14 have a maximum median recon-763 struction error of 29° (owl) and 47° (bust) due to the lighting 764 matrix being nearly singular. In addition to relaxing the cali-765 bration requirements (as full environment maps are not 766 needed for our technique), our learning-based technique 767 improves the median surface reconstruction accuracy by up 768 to 68 percent. 769

5.6 Analysis

We now analyze further our network, and in particular 771 explore the robustness of our network to departures from 772 the assumptions that were made in Section 5.2. *More analysis* 773 *is available in the supplementary material, available online* 774 *including results on a partially cloudy day.* 775

5.6.1 Camera Calibration Error

The impact on reconstruction performance when the north- 777 facing camera hypothesis is infringed was studied by rotat- 778 ing the real environment maps used to render the evalua- 779 tion dataset (Section 5.4), and show the results of this 780 experiment in Fig. 15 (left). The slight improvement around 781 5° west calibration error is due to the timestamps of our real 782 lighting dataset that are not perfectly aligned with the neural network expected timestamps. We observe that the 784 median reconstruction error increases of approximately 5° 785 per 10° error on camera calibration, showing that the network has some built-in robustness to these errors. 787

5.6.2 Number of Input Images

We also investigated normal estimation performance in func- 789 tion of the number of inputs *T* to the CNN (see Section 5.2.1). 790 Results ranging from a single input image (T = 1, effectively 791 performing shape-from-shading) to T = 16 input images 792 all uniformly taken from 9:00 to 16:00 are shown in 793 Fig. 15 (center). We observe an rapid improvement in perfor-794 mance from one to three images, which is coherent with PS 795 theory [1]. Performance continues to increase until T = 8, 796 probably because added constraints improves robustness to 797

770

776

788



Fig. 15. (Left) Median normal estimation error as box-percentile plots (see Fig. 13) in function of the camera deviation from north in degrees on our real lighting evaluation set. Positive means camera looking westward, negative means camera looking eastward. (Center) Normal estimation error as box-percentile plots on our evaluation dataset in function of the number of input images *T*. (Right) Ablative study on the number of pixels in input.

noise and non-diffuse materials. With T > 8, normal estimation error increases slightly. This could be due to an increase in the number of parameters to train in our model (the output tensor after concatenation is of dimension $14 \times 14 \times 32T$, thereby increasing the number of parameters in the second convolutional layer), making the model harder to train.

804 5.6.3 Patch Size

Fig. 15 (right) shows the impact of varying the patch size *P*. 805 To achieve this, we add an adaptive max pooling layer of 806 size 4×4 after the last residual block (see Fig. 10). Accuracy 807 increases with patch size until roughly 14 pixels, and then 808 decreases. We hypothesize that very small patch sizes do 809 not contain enough spatial context while too large patches 810 reveal macroscopic object features, which the neural net-811 work fails to recognize in the new shapes of the test set. 812

813 5.7 Limitations

The first limitation of the proposed method is that the cam-814 era is assumed to be pointing north. Although the network 815 shows some resilience to errors in camera calibration (see 816 Fig. 15), larger deviations from the assumed direction yield 817 degraded performance. One possible way to circumvent this 818 limitation would be to train direction-specific models and 819 select the right one by detecting the camera orientation. Fur-820 thermore, while our approach is robust to non-Lambertian 821 reflections, it assumes the scene to have a spatially-uniform 822 BRDF. This assumption is shared with recent techniques 823 like [18]. Fig. 16 shows the behavior with a spatially-varying 824 BRDF composed of a checkerboard pattern with small and 825 large squares. Unsurprisingly, the resulting normal maps 826



Fig. 16. Limitation of our approach. Our network is trained on spatially uniform BRDFs, so testing it on spatially-varying albedo maps increases the estimation error. (left) Spatially-uniform albedos results in low error, while checkerboard albedo maps with (center) small and (right) large patterns increase the error.

appear distorted since the constant albedo assumption is broken. One interesting direction for future work here would be to train a network on the *ratio* between pairs of images (e.g., as in [8]), which effectively cancels out the albedo.

6 **DISCUSSION**

This paper has presented a thorough analysis of outdoor PS ⁸³² under various illumination conditions captured over the ⁸³³ course of a single day. In this scenario, we have no control ⁸³⁴ over illumination, so existing methods for setting up opti-⁸³⁵ mal lighting [14], [15] cannot be applied. Through a data-⁸³⁶ driven analysis of the expected behavior of outdoor PS, we ⁸³⁷ reveal natural factors that distinguish good and unfavorable daylight conditions and identify mainly two different types ⁸³⁹ of working weather conditions: partially cloudy and clear ⁸⁴⁰ days. Our analysis shows that occlusion of the sun by ⁸⁴¹ clouds provides additional photometric cues that improve ⁸⁴² the accuracy of the surface reconstruction. Furthermore, ⁸⁴³ this improvement in conditioning can be observed in short ⁸⁴⁴ time intervals and varies in accord with surface orientation.

However, with a cloudless sky, outdoor PS becomes illconditioned (even in case of simple Lambertian reflectance) 847 and cannot be solved from photometric cues alone. To 848 address this issue, we augment the available photometric 849 cues with learned priors. As such, we present the first 850 method for single-day outdoor PS based on deep learning. 851 This new method is not limited to Lambertian objects and is 852 also robust to shadows and specular highlights. It signifiscantly outperforms previous work on a challenging evaluation dataset of virtual objects (lit by real sunny lighting 855 conditions) and yields successful surface reconstructions on real objects. 857

One exciting direction for future work is to leverage the 858 findings in this paper to design a unified approach for out-859 door PS under skies with any amount of cloud coverage. 860 For this, a properly-trained neural network could learn to 861 reconstruct the detailed surface of a large class of objects 862 observed under variable (but uncontrolled) natural outdoor 863 illumination. Designing such an approach will however 864 require care as simply training on clear days cannot reach 865 the same performance on other weather conditions (see the 866 supplementary material, available online). We hypothesize 867 that a hybrid approach, combining the photometric cues 868 from Section 4 with a deep CNN such as the one in Section 5 869 could be successful. A question that remains open to investigation is the adequate requirement in terms of lighting cal-871 ibration as to provide beneficial information during 872

reconstruction while also allowing for application in the wild. We believe the analysis presented in this paper sets

the stage for exciting future work.

876 **ACKNOWLEDGMENTS**

This research was supported by the NSERC Discovery GRANT RGPIN-2014-05314, the FRQ-NT New Researcher

879 Grant 2016-NC-189939, and the FRQ-NT REPARTI Strategic

880 Network. The authors also thank Nvidia for the donation of

the GPUs used in this research.

882 **REFERENCES**

- R. J. Woodham, "Photometric method for determining surface orientation from multiple images," *Optical Eng.*, vol. 19, no. 1, pp. 139–144, 1980.
- [2] N. G. Alldrin, T. Zickler, and D. J. Kriegman, "Photometric stereo with non-parametric and spatially-varying reflectance," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2008, pp. 1–8.
- [3] R. Basri, D. Jacobs, and I. Kemelmacher, "Photometric stereo with general, unknown lighting," *Int. J. Comput. Vis.*, vol. 72, no. 3, pp. 239–257, Jun. 2007.
- [4] C.-H. Hung, T.-P. Wu, Y. Matsushita, L. Xu, J. Jia, and C.-K. Tang,
 "Photometric stereo in the wild," in *Proc. IEEE Winter Conf. Appl. Comput. Vis.*, 2015, pp. 302–309.
- [5] J. Ackermann, F. Langguth, S. Fuhrmann, and M. Goesele,
 "Photometric stereo for outdoor webcams," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2012, pp. 262–269.
- A. Abrams, C. Hawley, and R. Pless, "Heliometric stereo: Shape from sun position," in *Proc. Eur. Conf. Comput. Vis.*, 2012, pp. 357–370.
- 900[7]E. Reinhard, G. Ward, S. Pattanaik, and P. Debevec, *High Dynamic Range Imaging*. Burlington, MA, USA: Morgan Kaufman, 2005.
- [8] L.-F. Yu, S.-K. Yeung, Y.-W. Tai, D. Terzopoulos, and T. F. Chan,
 "Outdoor photometric stereo," in *Proc. IEEE Int. Conf. Comput. Photography*, 2013, pp. 1–8.
- Photography, 2013, pp. 1–8.
 J.-F. Lalonde *et al.*, "The Laval HDR sky database," 2016. [Online].
 Available: http://www.hdrdb.com
- [10] B. Shi, Z. Mo, Z. Wu, D. Duan, S. K. Yeung, and P. Tan, "A bench-mark dataset and evaluation for non-lambertian and uncalibrated photometric stereo," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 2, pp. 271–284, Feb. 2019.
- [11] R. Onn and A. Bruckstein, "Integrability disambiguates surface recovery in two-image photometric stereo," *Int. J. Comput. Vis.*, vol. 5, no. 1, pp. 105–113, 1990.
- [12] C. Hernández, G. Vogiatzis, and R. Cipolla, "Overcoming shadows in 3-source photometric stereo," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 2, pp. 419–426, Feb. 2011.
- [13] F. Shen, K. Sunkavalli, N. Bonneel, S. Rusinkiewicz, H. Pfister, and
 X. Tong, "Time-lapse photometric stereo and applications,"
 Comput. Graphics Forum, vol. 33, no. 7, pp. 359–367, 2014.
- [14] O. Drbohlav and M. Chantler, "On optimal light configurations in photometric stereo," in *Proc. 10th IEEE Int. Conf. Comput. Vis.*, 2005, pp. 1707–1712.
- M. Klaudiny and A. Hilton, "Error analysis of photometric stereo
 with colour lights," *Pattern Recognit. Lett.*, 2014. [Online]. Avail able: http://www.sciencedirect.com/science/article/pii/S0167
 865513005114
- P. Debevec, "Rendering synthetic objects into real scenes: Bridging traditional and image-based graphics with global illumination and high dynamic range photography," in *Proc. 25th Annu. Conf. Comput. Graphics Interactive Tech.*, 1998, pp. 189–198.
- [17] B. Shi, K. Inose, Y. Matsushita, P. Tan, S.-K. Yeung, and K. Ikeuchi,
 "Photometric stereo using Internet images," in *Proc. Int. Conf. 3D* Vis., 2014, pp. 361–368.
- [18] Z. Mo, B. Shi, F. Lu, S.-K. Yeung, and Y. Matsushita, "Uncalibrated photometric stereo under natural illumination," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 2936–2945.
- [19] K. Inose, S. Shimizu, R. Kawakami, Y. Mukaigawa, and K. Ikeuchi,
 "Refining outdoor photometric stereo based on sky model," *Inf. Media Tech.*, vol. 8, no. 4, pp. 1095–1099, Dec. 2013.
- [20] J. Jung, J.-Y. Lee, and I. S. Kweon, "One-day outdoor photometric stereo using skylight estimation," Int. J. Comput. Vis., vol. 127, pp. 1126–1142, 2019.

- [21] G. Oxholm and K. Nishino, "Shape and reflectance from natural 943 illumination," in Proc. Eur. Conf. Comput. Vis., 2012, pp. 528–541. 944
- M. K. Johnson and E. H. Adelson, "Shape estimation in natural 945 illumination," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 946 2011, pp. 2553–2560.
- [23] J. T. Barron and J. Malik, "Shape, illumination, and reflectance 948 from shading," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, 949 no. 8, pp. 1670–1687, Aug. 2015. 950
- [24] D. Eigen and R. Fergus, "Predicting depth, surface normals and 951 semantic labels with a common multi-scale convolutional 952 architecture," in Proc. IEEE Int. Conf. Comput. Vis., 2015, 953 pp. 2650–2658. 954
- pp. 2650–2658.
 [25] Z. Shu, E. Yumer, S. Hadap, K. Sunkavalli, E. Shechtman, and 955
 D. Samaras, "Neural face editing with intrinsic image disen-956 tangling," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2017, 957
 pp. 5444–5453.
- [26] J. Wu, Y. Wang, T. Xue, X. Sun, W. T. Freeman, and J. B. Tenenbaum, 959
 "Marrnet: 3D shape reconstruction via 2.5D sketches," in *Proc.* 960
 Advances Neural Inf. Process. Syst., 2017, pp. 540–550. 961
- [27] J.-F. Lalonde and I. Matthews, "Lighting estimation in outdoor 962 image collections," in Proc. IEEE Int. Conf. 3D Vis., 2014, pp. 131–138. 963
- Y. Hold-Geoffroy, J. Zhang, P. F. U. Gotardo, and J.-F. Lalonde, 964
 "What is a good day for outdoor photometric stereo?" in *Proc. Int.* 965 *Conf. Comput. Photography*, 2015, pp. 1–9. 966
- Y. Hold-Geoffroy, J. Zhang, P. F. U. Gotardo, and J.-F. Lalonde, 967
 "x-hour outdoor photometric stereo," in *Proc. Int. Conf. 3D Vis.*, 968
 2015, pp. 28–36.
- [30] Y. Yu and W. A. P. Smith, "PVNN: A neural network library for 970 photometric vision," in Proc. IEEE Int. Conf. Comput. Vis., 2017, 971 pp. 526–535. 972
- [31] H. Santo, M. Samejima, Y. Sugano, B. Shi, and Y. Matsushita, 973
 "Deep photometric stereo network," in *Proc. IEEE Int. Conf.* 974
 Comput. Vis., 2017, pp. 501–509. 975
- [32] T. Taniai and T. Maehara, "Neural photometric stereo reconstruction for general reflectance surfaces," 2018, arXiv: 1802.10328. 977
- [33] J. Sun, M. Smith, L. Smith, and A. Farooq, "Examining the uncertainty of the recovered surface normal in three light photometric stereo," *Image Vis. Comput.*, vol. 25, no. 7, pp. 1073–1079, 2007.
- [34] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical* 981 *Learning: Data Mining, Inference, and Prediction*. NY, USA: Springer, 982 2009.
- [35] D. F. Dementhon and L. S. Davis, "Model-based object pose in 25 984 lines of code," Int. J. Comput. Vis., vol. 15, no. 1–2, pp. 123–141, 1995. 985
- [36] L. Hošek and A. Wilkie, "An analytic model for full spectral skydome radiance," ACM Trans. Graphics, vol. 31, no. 4, pp. 1–9, 2012. 987
- [37] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning 988 for image recognition," in Proc. IEEE Conf. Comput. Vis. Pattern 989 Recognit., 2016, pp. 770–778. 990
- [38] D.-A. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (elus)," in 992 Proc. Int. Conf. Learn. Representations, 2016. 993
- [39] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image 994 translation with conditional adversarial networks," in *Proc. IEEE* 995 *Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 1125–1134. 996
- [40] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep 997 network training by reducing internal covariate shift," in *Proc. Int.* 998 *Conf. Mach. Learn.*, 2015, pp. 448–456. 999
 [41] D. Kingma and J. Ba, "Adam: A method for stochastic opti-100
- [41] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," in Proc Int. Conf. Learn. Representations, 2015. 1001
- B. Curless and M. Levoy, "A volumetric method for building complex models from range images," in *Proc. 23rd Annu. Conf. Comput.* 1003 *Graphics Interactive Tech.*, 1996, pp. 303–312.
- [43] B. Walter, S. Marschner, H. Li, and K. Torrance, "Microfacet 1005 models for refraction through rough surfaces," in *Proc. 18th* 1006 *Eurographics Conf. Rendering Tech.*, 2007, pp. 195–206. 1007
- [44] P. Bretagnon and G. Francou, "Planetary theories in rectangular 1008 and spherical variables-vsop 87 solutions," Astronomy Astrophysics, vol. 202, pp. 309–315, 1988.
 1010
- [45] J. Stumpfel *et al.*, "Direct HDR capture of the sun and sky," in *Proc.* 1011 *AFRIGRAPH*, 2004, pp. 145–149.
 1012
- [46] W. W. Esty and J. D. Banfield, "The box-percentile plot," J. Statisti- 1013 cal Softw., vol. 8, pp. 1–14, 2003. [Online]. Available: http://www. 1014 jstatsoft.org/v08/i17/paper 1015

HOLD-GEOFFROY ET AL.: SINGLE DAY OUTDOOR PHOTOMETRIC STEREO

1025

1026

1027

1028

1029

1030 1031

1032

1033

1034

1035



Yannick Hold-Geoffroy received the PhD degree in electrical engineering with a mention on the dean's honor roll from Université Laval, Canada, in 2018. He is currently a research engineer with Adobe in San José. He was awarded the CIPPRS Doctoral Dissertation Award in 2019. His research interests include understanding of natural images through machine learning and lighting modelling and estimation.



Jean-François Lalonde received the MS and 1042 PhD degrees from Carnegie Mellon, in 2006 and 1043 2011, respectively. He is currently an associate 1044 professor with the ECE Department, Université 1045 Laval, Canada. Previously, he was a postdoctoral 1046 associate at Disney Research, Pittsburgh. His PhD 1047 thesis won the 2010-11 CMU School of Computer 1048 Science Distinguished Dissertation Award. His 1049 research interests include computer vision and 1050 deep learning, with a particular focus on lighting estimation, 3D reconstruction, tracking, and augmented 1052 reality. 1053

▷ For more information on this or any other computing topic, 1054 please visit our Digital Library at www.computer.org/csdl.

13

Paulo Gotardo received the BSc and MSc degrees in informatics from the Federal University of Parana, Brazil, and the PhD degree in electrical and computer engineering from the Ohio State University. He is currently a senior research scientist with Disney Research Studios in Zurich. His research on computer vision, graphics and machine learning focuses on modeling and capturing 3D geometry, motion, appearance and illumination in dynamic scenes, with application in building digital humans (mov-

ies, games, AR, and VR). Before joining Disney Research, he was a
research associate with the Advanced Computing Center for the Arts
and Design and a postdoctoral researcher with the Computational Biology and Cognitive Science Lab, both at Ohio State. He was also an
associate research scientist with disney research pittsburgh on the Carnegie Mellon University campus.