

Integral imaging techniques for flexible sensing through image-based reprojection

JOSE M. SOTOCA,^{1,*} PEDRO LATORRE-CARMONA,¹ FILIBERTO PLA,¹ XIN SHEN,² SATORU KOMATSU,² AND BAHRAM JAVIDI²

¹Institute of New Imaging Technologies, Universitat Jaume I, 12071 Castellón de la Plana, Spain

²Electrical and Computer Engineering Department, University of Connecticut, Storrs, Connecticut 06269-4157, USA

*Corresponding author: sotoca@uji.es

Received 9 May 2017; revised 7 August 2017; accepted 15 August 2017; posted 16 August 2017 (Doc. ID 295585); published 0 MONTH 0000

In this work, a 3D reconstruction approach for flexible sensing inspired by integral imaging techniques is proposed. This method allows the application of different integral imaging techniques, such as generating a depth map or the reconstruction of images on a certain 3D plane of the scene that were taken with a set of cameras located at unknown and arbitrary positions and orientations. By means of a photo-consistency measure proposed in this work, *all-in-focus* images can also be generated by projecting the points of the 3D plane into the sensor planes of the cameras and thereby capturing the associated RGB values. The proposed method obtains consistent results in real scenes with different surfaces of objects as well as changes in texture and lighting. © 2017 Optical Society of America

OCIS codes: (100.6890) Three-dimensional image processing; (110.3010) Image reconstruction techniques; (110.6880) Three-dimensional image acquisition.

<https://doi.org/10.1364/JOSAA.99.099999>

1. INTRODUCTION

As opposed to traditional two-dimensional (2D) imaging techniques, three-dimensional (3D) imaging technologies can potentially capture the 3D structure, range, and texture information of the different objects in a scene. Additionally, 3D imaging technologies are more robust to partial scene occlusion. There are many 3D imaging technologies, such as holography and related interferometry techniques [1], stereoscopy [2], pattern illumination techniques [3], LADAR [4], and time-of-flight techniques [5].

Multi-perspective imaging obtains 3D scene information by recording conventional 2D incoherent images from multiple views. Because standard 2D images are used, multi-perspective 3D imaging systems can be built using a single inexpensive camera with a lenslet array or an array of inexpensive sensors. However, thanks to the advances in optoelectronic sensors such as CMOS and CCDs, display devices such as LCDs, and commercially available digital computers, integral imaging is a very active area of research nowadays.

Integral imaging can be considered a class of multi-view imaging acquisition and display technology [6]. It has been applied in fields like visualization [7], target recognition and ranging [8], 3D photon-counting imaging [9,10], 3D imaging for objects under occlusions or in a scattering medium [11], 3D underwater imaging [12], biological or medical imaging [13],

integral microscopy [14], and others [15], to cite just a few examples.

Integral imaging performs well under ambient or incoherent light, which compares favorably in relation to other sensing techniques, such as holography, LADAR, or structured light, that make use of an active illumination system. It also has specific benefits over 2D imaging as well as stereo imaging. For 3D visualization purposes in integral imaging, the microlenses produce differences in the light density within the space in front of the observer. Thus, there is a real reconstruction of the light structure produced by the original 3D scene. In lenslet-based integral imaging systems, the achievable resolution is limited by the size of the lenslet and the number of pixels allocated to each lenslet. In essence, the resolution of each elemental image (EI) is limited by three parameters: the pixel size, the lenslet point spread function, and the lenslet depth of focus [7,16]. In addition, aberrations and diffraction are significant because the size of the lenslet is relatively small. In contrast to the lenslet-based systems, integral imaging can be performed either in a synthetic aperture mode or with an array of high-resolution imaging sensors. Each perspective image can be recorded by a full-size CCD or CMOS sensor of several megapixels. This approach may be considered *synthetic aperture integral imaging* (SAII) [17]. SAII enables larger fields of view (FOVs) to be obtained with high resolution 2D images because each 2D image makes full use of the detector array and the optical aperture.

Traditionally, SAI consists of a setup formed by a camera array located on a planar surface. This configuration greatly limits the application of SAI in other situations where an arbitrary arrangement of the cameras is necessary. In Ref. [18], the authors propose the use of a lenslet-based integral imaging system that has an array of lenslets embedded in an elastic scaffold, integrating it into a flexible optoelectronic detector array in an arbitrary non-planar configuration. Recent advances in mechanics and material properties of conventional rigid wafer-based technologies, but with the ability to be stretched and deformed into arbitrary shapes, allow active components to be connected to create new engineering options in imaging devices, where the geometry of the detector array can be optimized together with the lens configuration [19]. The most promising initial possibilities for application are in surveillance, night vision, endoscopy, and retinal implants, or as active components on the eye to enhance vision.

Moreover, authors in Ref. [20] present a 3D integral image acquisition and reconstruction technique with unknown sensor positions and orientations placed on a flexible surface that increases the field of view of the 3D imaging system. In addition, the proposed estimation algorithm assumes that the relative pose of the first two cameras is known. This may not be very convenient if we want to carry out experiments in real-world scenarios without any constraints. Another problem that arises when seeking to solve 3D reconstruction with sensors on a flexible surface is how to obtain a criterion that is robust in this type of problem with an arbitrary camera arrangement. In Ref. [21], a methodology is developed to build a depth map of the scene using a minimum variance approach. Depth estimation accuracy will degrade when object surfaces do not satisfy the Lambertian assumption and requires a precise photometric calibration of the cameras, such as in the presence of partial occlusions or when concave surfaces exist. To address this problem, several proposals for multi-view photo-consistency measures have been developed, such as voxel coloring [22], space carving [23], standard deviation based on an adaptive threshold [24], and voting strategies [25], and in some deformable surface methods [26]. Similarly, the variational formulation relies on square intensity differences [27] or modeling the intensity deviations from brightness constancy by a multivariate Gaussian [28]. Other photo-consistency measures are based on the assumption that a comparison can be made between pairs of images used in stereo, such as normalized cross-correlation (NCC), the sum of squared differences (SSD), mutual information-based measures, and others [29]. Nevertheless, this does not remove any of the severe limitations of the Lambertian assumption.

Integral imaging offers a series of advantages in relation to other 3D imaging techniques. Three of them are: (a) its capability to reconstruct a scene on planes at a constant depth, where only the objects that are at that distance from the camera array are *in focus*; (b) the creation of an *all-in-focus* image from the stack of depth planes; and (c) the ability to infer a depth map of the scene.

The main contribution of this work is oriented toward providing a technique to adapt the reconstruction methodology applied in integral imaging for a *flexible sensing* configuration

and show that these same features (i.e., focus on a given depth, creation of an *all-in-focus* image, estimation of the scene depth, etc.) can be obtained in a *flexible sensing* configuration. To that end, a precise calibration of the system is used based on [30], which does not need knowledge of any intrinsic or extrinsic parameter of the cameras setup.

To show the feasibility and accuracy of the proposed 3D plane reconstruction by reprojection, we analyze the problem to obtain a photo-consistency criterion introduced in Section 3 for flexible sensing setups. Thus, we apply the approach proposed in Ref. [31] for light field displays, consisting in a defocusing strategy to deal with spatial information surrounding a pixel. Furthermore, a comparison of this photo-consistency measure will be made with the method based on minimum variance that has been widely used in integral imaging.

Although in Ref. [31], an occlusion method for light fields is also proposed, this is not applicable to the case of an arbitrary flexible sensing setup due to the amount of disparity among the elemental images, which makes the occlusion problem worse than in usual integral imaging setups. In this sense, we have chosen an alternative occlusion method proposed in Ref. [32] and explained in Section 4.

The rest of the paper is organized as follows. Section 2 provides a brief explanation of how the calibration process has been solved in an arbitrary cameras setup. Section 3 explains the methodology proposed in this paper for robust depth estimation. Section 4 describes the creation of the depth map and the *all-in-focus* image estimation of the scene. Section 5 offers the results obtained by applying the techniques proposed here in real scenes and also discusses several aspects. Finally, several conclusions are given in Section 6.

2. MULTI-CAMERA SELF-CALIBRATION

Important advances have been recently made in the reconstruction of 3D scenes from multiple views. In this sense, the review by Ref. [33] and the associated Middlebury evaluation framework represented a milestone after which a lot of research has been conducted focusing on the multi-view reconstruction of objects taken under strictly controlled imaging conditions. However, most of these algorithms are not directly suited to large-scale outdoor scenes.

In a multi-view camera acquisition system, we also need a calibration algorithm that is sufficiently precise to be used for integral imaging techniques that may be able to perform well in outdoor scenes. In this section, we describe a calibration method and camera location for the case where the cameras have an arbitrary pose. The method is based on the work proposed in Ref. [30] for m -views using bundle adjustment for a projective reconstruction.

Consider the case where a set of n 3D points $\mathbf{X}_j = [X_j, Y_j, Z_j, 1]^T, j = 1, \dots, n$ are viewed by a set of m cameras with projection matrices P^i . Denote by $\mathbf{x}_j^i = [x_j^i, y_j^i, 1]^T$ the coordinates of the j -th point as seen in the i -th camera. Our goal is to solve the reconstruction problem where, given a set of image coordinates \mathbf{x}_j^i , we aim to find the set of camera matrices P^i and their correspondence points \mathbf{X}_j in the scene such that

$$\mathbf{x}_j^i = P^i \mathbf{X}_j. \quad (1)$$

183 If the image correspondence measurements has a high number
 184 of uncertainties, then Eq. (1) will not be satisfied exactly.
 185 **2** Thus, we wish to estimate the projection matrices \hat{P}^i and the
 186 3D points $\hat{\mathbf{X}}_j$ that project exactly onto image points $\hat{\mathbf{x}}_j^i$ as
 187 $\hat{\mathbf{x}}_j^i = \hat{P}^i \hat{\mathbf{X}}_j$. Likewise, we also seek to minimize the image distance
 188 between the reprojected points and detected (measured)
 189 image points \mathbf{x}_j^i for every view in which the 3D point is seen.
 190 This approximation (which minimizes the reprojection error) is
 191 defined as *bundle adjustment* [34].

192 Equation (1), representing the projective mapping, can be
 193 interpreted as true only up to a constant factor. Writing this
 194 constant factor explicitly, we have

$$\lambda_j^i \mathbf{x}_j^i = P^i \mathbf{X}_j. \quad (2)$$

195 Thus, the goal of the calibration is to estimate the scales λ_j^i and
 196 the camera projection matrices P^i . The weighting factors λ_j^i are
 197 called the *projective depths* of the points.

198 For the estimation of λ_j^i we have used Sturm and Triggs'
 199 method, exploiting epipolar geometry to obtain these projective
 200 depths [35]. In relation to the two alternative methodologies
 201 proposed by the authors, the solution based on a central image
 202 is more appropriate for wide baseline stereo, and it is the one
 203 we used in this section (see Martinec and Pajdla [36] for more
 204 details).

205 Provided that each point is visible in every view, we can put
 206 all the points and camera projections into the W_s matrix:

$$W_s = \begin{bmatrix} \lambda_1^1 \mathbf{x}_1^1 & \dots & \lambda_n^1 \mathbf{x}_n^1 \\ \vdots & \vdots & \vdots \\ \lambda_1^m \mathbf{x}_1^m & \dots & \lambda_n^m \mathbf{x}_n^m \end{bmatrix} = \begin{bmatrix} P^1 \\ \vdots \\ P^m \end{bmatrix} [\mathbf{X}_1 \dots \mathbf{X}_n], \quad (3)$$

207 where W_s is called the *scaled measurement matrix*, $P =$
 208 $[P^1 \dots P^m]^T$ and $X = [\mathbf{X}_1 \dots \mathbf{X}_n]$. P and X are referred to as
 209 the *projective motion* and the *projective shape*, respectively [30].

210 If we collect enough noiseless points (x_j^i, y_j^i) and the scales λ_j^i
 211 are known, then W_s can be factored into \hat{P} and X [35]. The
 212 factorization of Eq. (3) retrieves the motion and shape through
 213 a 4×4 projective transformation H :

$$W_s = PX = PHH^{-1}X = \hat{P}\hat{X}, \quad (4)$$

214 where $\hat{P} = PH$ and $\hat{X} = H^{-1}X$. The self-calibration process
 215 computes a matrix H , such that \hat{P} and \hat{X} become Euclidean.
 216 This process is sometimes called *Euclidean stratification* [30,37].
 217 The matrix H can be solved by imposing certain geometrical
 218 constraints. The most general constraint is the assumption that
 219 some internal parameters of the cameras are the same.

220 In Ref. [36], projective reconstruction by factorization is
 221 applied, handling perspective views and occlusions jointly.
 222 The factorization algorithm also provides an optimal method
 223 for computing the new image points when they are not visible
 224 from all the cameras. In addition, the method proposed in
 225 Ref. [30] fills the missing points (those with unknown depths).
 226 This is implemented in two steps: first, triangulation to find the
 227 pre-image X , and then the reprojection as PX to generate its
 228 image in all views. In practice, triangulation and reprojection
 229 provide a method of "filling in" points that are missed during
 230 multiple view matching.

231 Another aspect to be considered is that lenses with short focal
 232 lengths are often used in immersive environments to guarantee

sufficient field of view. However, such lenses have significant
 nonlinear distortion, which has to be corrected for precise 3D
 computation. Therefore, a distortion model is applied to assess
 the radial and tangential distortion, aiming at eliminating these
 distortion effects of the lenses in the elemental images obtained
 by the cameras during the calibration process.

3. PHOTO-CONSISTENCY RECONSTRUCTION BY IMAGE-BASED REPROJECTION

In integral imaging, an optical display or computational
 reconstruction method can be used to visualize a 3D scene.
 In the computational reconstruction approach, the elemental
 images obtained during the acquisition stage are projected onto
 the image plane at an arbitrary distance through a real pinhole
 or lens. Because a 3D object can be viewed as the combination
 of multiple depth images, 3D information can be observed and
 analyzed by generating a series of depth images.

For a flexible sensing setup, as is our case, we adapt the com-
 putational reconstruction used in integral imaging for a regular
 array of sensors to the case of a non-uniformly distributed flex-
 ible sensing integral imaging system, where the camera setup is
 not placed on a flat surface with known positions in a regular
 grid. Thus, an alternative strategy is to sweep a set of planes
 through the scene with respect to a reference camera [see
 Fig. 1(a)]. This is known as the *plane sweep* algorithm in the
 computer vision literature [38,39]. Sweeping to a depth D
 through a series of disparity hypotheses corresponds to mapping
 each input image into the reference camera defining the disparity
 space through a series of homographic transformations [38].

We have the projection matrices of the different cameras
 obtained by the calibration method explained in the previous
 section. Therefore, the approach presented here is aimed at
 achieving a depth reconstruction of the objects that are ob-
 served from this reference camera c , which we call the *central*
camera, and whose projection matrix is defined by $P^c: \mathbb{R}^3 \rightarrow \mathbb{R}^2$.
 We denote $I^c: \Omega_c \subset \mathbb{R}^2 \rightarrow \mathbb{R}^d$ as the intensity of the image
 acquired by camera c in the set of pixels Ω_c . In practice, the
 parameter d defines the information stored in the pixels by tak-
 ing the value $d = 1$ for grayscale images and $d = 3$ for RGB
 images. Thus,

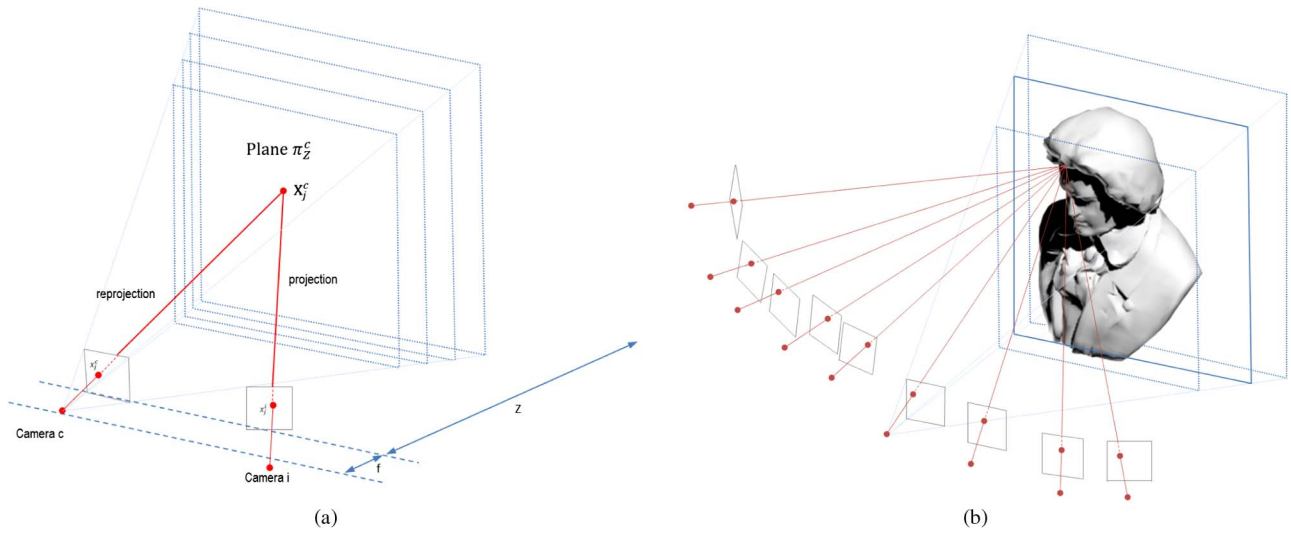
$$\mathbf{X}_j^c = P_{\pi_Z^c}^{-1} \lambda_Z^c \mathbf{x}_j^c. \quad (5)$$

Let us consider that from the camera c we want to reproject
 the set of pixels Ω^c onto a 3D plane called *Plane* π_Z^c , which is
 located at a distance Z with respect to its optical center.
 Therefore, let us define a reprojection from the camera c onto
 the plane by $P_{\pi_Z^c}^{-1}: P^c(\Omega_c) \rightarrow \pi_Z^c$.

During the calibration process, each detected 3D point in
 the scene [see Eq. (2)] has a scaling factor λ_j^i that is different,
 and it depends on the depth in relation to the camera. To gener-
 ate a 3D plane at a certain depth, we must only use a constant
 scale factor $\lambda_Z^c = Z$. Depending on the applied factor in
 $[Z_{\min}, Z_{\max}]$ on the image pixels, we can generate 3D points
 \mathbf{X}_j^c on planes π_Z^c located at different depth ranges, as shown
 in Fig. 1(a).

Every time we reproject the pixels of the image I^c to a depth
 level Z , these 3D points can be seen by the rest of the cameras,
 and, therefore, their positions on their respective images can be

233
 234
 235
 236
 237
 238
 239
 240
 241
 242
 243
 244
 245
 246
 247
 248
 249
 250
 251
 252
 253
 254
 255
 256
 257
 258
 259
 260
 261
 262
 263
 264
 265
 266
 267
 268
 269
 270
 271
 272
 273
 274
 275
 276
 277
 278
 279
 280
 281
 282
 283
 284
 285
 286
 287



F1:1 **Fig. 1.** (a) Reprojection and projection operations with respect to a point \mathbf{X}_j^c on the plane π_z^c . (b) Camera setup, with arbitrary distribution
F1:2 observing a point in the scene.

288 estimated. Thus, the value of the image observed by camera i
289 via reprojecting of the camera c will be expressed by

$$I^i \circ P^i \circ P_{\pi_z^c}^{-1} \circ P^c(\Omega_c) \rightarrow \mathbb{R}^d. \quad (6)$$

290 When the plane is at the depth corresponding to the distance
291 that the object is situated with respect to the optical center of
292 the *center camera* and in the absence of occlusions, we can con-
293 sider that all cameras will observe the same image value [see
294 Fig. 1(b)], which is the principle of photo-consistency.

295 A. Photo-Consistency Measure

296 As indicated in the introduction, an accurate scene depth
297 estimation may be degraded by the shape of the objects or
298 by the intersection of objects with others seen by the different
299 cameras. Hence, the matching process between the different
300 views must handle projective distortion and partial occlusions.
301 The use of local as well as global image intensity information
302 can be exploited to improve the robustness to changes in ap-
303 pearance, without taking into account any approximation of
304 shape, motion, or visibility.

305 An example of occlusion due to a convex shape can be seen
306 in Fig. 2. We show a scheme with four cameras C_1 to C_4 , pro-
307 ducing four images of an object with intensities I^1 to I^4 . Each
308 3D point \mathbf{X} projects on the positions \mathbf{x}_1 to \mathbf{x}_4 of the images in
309 the cameras. Model (object) point \mathbf{X} projects to \mathbf{x}_2 and \mathbf{x}_3
310 with intensities I^2 and I^3 but not in \mathbf{x}_1 and \mathbf{x}_4 . Thus, the intensities
311 of I^1 and I^4 are not equal to I^2 and I^3 . Needless to say, this
312 is just one of the ways in which occlusion occurs, and other
313 combinations can be produced.

314 The previous example does not satisfy the conditions of a
315 Lambertian lighting model, where image intensity or color
316 per pixel would be independent of the camera viewpoint.
317 Therefore, the image intensities at pixels \mathbf{x}_1 to \mathbf{x}_4 should be
318 identical apart from image noise and differences in the camera
319 responses. Let a set of optical images with intensity values be
320 (I^1, \dots, I^m) . Thus, we can project each 3D point of the object
321 \mathbf{X}_j onto the corresponding pixel \mathbf{x}_j^i for each camera i . Then, the

322 arithmetic mean associated to the pixel values of the images
323 corresponding to an object point would be given by:

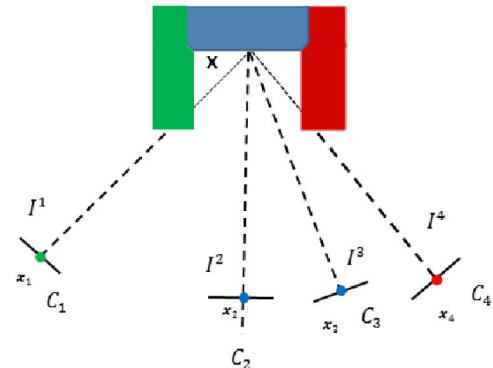
$$\bar{I}_j = \frac{1}{m} \sum_{i=1}^m I^i(\mathbf{x}_j^i). \quad (7)$$

324 The variance between image intensities and the mean is
325 defined as

$$V_j^2 = \frac{1}{m} \sum_{i=1}^m (I^i(\mathbf{x}_j^i) - \bar{I}_j)^2. \quad (8)$$

326 The variance criterion was one of the first photo-consistency
327 criteria proposed and is also one of the most widely accepted.

328 Because of the errors in the photo-consistency estimation
329 introduced by occlusions, and taking into account that this fact
330 worsens in an arbitrary flexible sensing setup, the previous
331 arithmetic mean [see Eq. (7)] is computed by applying a visi-
332 bility criterion for each pixel $O_{occ}^i(\mathbf{x}_j^i, \pi_z^c)$ that takes values 0 or
333 1, considering two conditions. The first condition establishes
334 that for those 3D points in the plane π_z^c that are not visible
335 for the other cameras, the visibility criterion of pixels takes the



F2:1 **Fig. 2.** Projection with occlusions. All the cameras do not see the
F2:2 same point in the scene.

value 0, otherwise it takes 1. As a second condition, we use the asymmetrical occlusion model of Wei and Quan [32] to evaluate the visibility of pixels $O_{occ}^i(\mathbf{x}_j^i, \pi_Z^c)$ for each level of Z . It is defined as being 0 if there exists another pixel \mathbf{p}_j^i which projects onto the same point in camera i as pixel \mathbf{x}_j^i and for which the projected depth is less than that of \mathbf{x}_j^i , otherwise it is 1. Therefore, the arithmetic mean associated to the pixel of the *central camera* and its reprojection on the plane π_Z^c is given by:

$$\bar{I}_{\pi_Z^c(j)} = \frac{\sum_{i=1}^m I^i(\mathbf{x}_j^i) \cdot O_{occ}^i(\mathbf{x}_j^i, \pi_Z^c)}{\sum_{i=1}^m O_{occ}^i(\mathbf{x}_j^i, \pi_Z^c)}. \quad (9)$$

In addition, the variance between image intensities would be:

$$V_{\pi_Z^c(j)}^2 = \frac{\sum_{i=1}^m (I^i(\mathbf{x}_j^i) - \bar{I}_{\pi_Z^c(j)})^2 \cdot O_{occ}^i(\mathbf{x}_j^i, \pi_Z^c)}{\sum_{i=1}^m O_{occ}^i(\mathbf{x}_j^i, \pi_Z^c)}. \quad (10)$$

346 Algorithm 1: Depth map scene and *all-in-focus* image

348 1: **Procedure** Flexible Sensing Integral Imaging by reprojection of 3D
 349 planes
 350 2: **Input:**
 351 3: Ω^c : set of pixels of *central camera*
 352 4: P^i : set of camera projection matrices
 353 5: I^i : set of images captured by the cameras
 354 6: $zstep$: distance in depth between two planes π_Z^c
 355 7: **Output:**
 356 8: $Z_a = Z_b = Z_{min}$
 357 9: **while** $Z_b \leq Z_{max}$ **do** \triangleright work with two depths Z_a and Z_b
 358 10: $Z_b = Z_b + zstep$
 359 11: \forall pixel \mathbf{x}_j^c in Ω^c estimate 3D points \mathbf{X}_{j,Z_a}^c in depth scene
 360 by reprojection
 361 12: \forall pixel \mathbf{x}_j^c in Ω^c estimate 3D points \mathbf{X}_{j,Z_b}^c in plane $\pi_{Z_b}^c$
 362 by reprojection
 363 13: $i = 1$
 364 14: **while** $i \leq m$ **do** $\triangleright m$ is the number of cameras
 365 15: \forall 3D point calculated, project $\mathbf{x}_{j,Z_a}^i = P^i \mathbf{X}_{j,Z_a}^c$
 366 16: \forall 3D point in plane $\pi_{Z_b}^c$, project $\mathbf{x}_{j,Z_b}^i = P^i \mathbf{X}_{j,Z_b}^c$
 367 17: \forall pixel \mathbf{x}_j^i store the intensities $I^i(\mathbf{x}_{j,Z_a}^i)$, $I^i(\mathbf{x}_{j,Z_b}^i)$.
 368 18: Thus, store the visibility $O_{occ}^i(\mathbf{x}_{j,Z_a}^i)$, $O_{occ}^i(\mathbf{x}_{j,Z_b}^i)$
 369 19: $i = i + 1$
 370 20: **return**
 371 21: \forall pixel \mathbf{x}_j^c estimate $\bar{I}_{\pi_{Z_a}^c(j)}$, $\bar{I}_{\pi_{Z_b}^c(j)}$, $V_{\pi_{Z_a}^c(j)}^2$, $V_{\pi_{Z_b}^c(j)}^2$
 372 22: **If** $Z_b = Z_{min} + zstep$
 373 23: \forall pixel \mathbf{x}_j^c estimate $\widehat{Photo}_{\pi_{Z_a}^c}(\mathbf{x}_j^c)$
 374 24: **EndIf**
 375 25: \forall pixel \mathbf{x}_j^c estimate $\widehat{Photo}_{\pi_{Z_b}^c}(\mathbf{x}_j^c)$
 376 26: **If** $\widehat{Photo}_{\pi_{Z_b}^c}(\mathbf{x}_j^c) < \widehat{Photo}_{\pi_{Z_a}^c}(\mathbf{x}_j^c)$
 377 27: $\widehat{Photo}_{\pi_{Z_a}^c}(\mathbf{x}_j^c) = \widehat{Photo}_{\pi_{Z_b}^c}(\mathbf{x}_j^c)$ and $Z_{step}(\mathbf{x}_j^c) = Z_b$
 378 28: **Else**
 379 29: $Z_{step}(\mathbf{x}_j^c) = Z_a$
 380 30: **EndIf**
 381 31: **If** $Z_b \leq Z_{max}$
 382 32: $Z_a \leftarrow Z_{step}$
 383 33: **EndIf**
 384 34: **return** $\bar{I}_{\pi_{Z_a}^c}, Z_a$

385 A drawback in this strategy is that it is usually applied as
 386 a per-pixel photo-consistency measure. To give more ro-
 387 bustness to noise and to be able to deal with realistic imaging

conditions, pixel neighborhood imaging information should be
 388 incorporated.

389 Authors in Ref. [31] estimate the depth of a scene by com-
 390 bining a *defocus* and a *correspondence* measure. However, they
 391 apply it to light fields where the object disparity in the elemental
 392 images is small. That is not our case. On the other hand, the
 393 defocus measure allows an optimal contrast in a certain region
 394 of the image to be obtained, but occlusions and lighting
 395 changes may easily affect the measurement accuracy. The patch
 396 size may also affect the measure sensitivity because the defocus
 397 measure may exceed the patch size. Correspondence measurement
 398 allows depth to be estimated using photo-consistency, and it has
 399 been widely used in stereo problems. In this case, a statistical mea-
 400 sure is usually applied to resolve matching ambiguities.

401 In our approach, we propose a photo-consistency measure
 402 where the first term (correspondence term) defines an initial
 403 cost function equal to the square root of the variance $V_{\pi_Z^c(j)}^2$
 404 [see Eq. (10)] in the plane π_Z^c for each point [see Fig. 1(a)].
 405 The second term (defocus term) acts locally and involves the
 406 reconstructed mean image intensities defined as $\bar{I}_{\pi_Z^c(j)}$ in rela-
 407 tion to the intensities I^c of the *central camera*. Thus, given a 3D
 408 point \mathbf{X}_j^c reprojected from pixel \mathbf{x}_j^c of the *central camera* in the
 409 plane π_Z^c , therefore,
 410

$$\text{Photo}_{\pi_Z^c}(\mathbf{x}_j^c) = \left\{ \sqrt{V_{\pi_Z^c(j)}^2} + |\bar{I}_{\pi_Z^c(j)} - I^c(\mathbf{x}_j^c)| \right\}. \quad (11)$$

411 Moreover, we add neighborhood imaging information by
 412 applying a bilfiltering technique with a spatial mean and a zero
 413 mean Gaussian kernel function G_s on the image intensity
 414 differences, centered on the current pixel around a window
 415 \mathcal{W} defined as

$$\widehat{\text{Photo}}_{\pi_Z^c}(\mathbf{x}_j^c) = \frac{\sum_{\mathbf{p}_i^c \in \mathcal{W}} \text{Photo}_{\pi_Z^c}(\mathbf{p}_i^c) \cdot G_s(|I^c(\mathbf{p}_i^c) - I^c(\mathbf{x}_j^c)|)}{\sum_{\mathbf{p}_i^c \in \mathcal{W}} G_s(|I^c(\mathbf{p}_i^c) - I^c(\mathbf{x}_j^c)|)}. \quad (12)$$

416 The idea of using color differences as a range filter to esti-
 417 mate the photo-consistency value is based on the observation
 418 that whenever a change of depth edge appears, a color change
 419 usually occurs between background objects with respect to fore-
 420 ground objects. This can be useful for comparing neighbor-
 421 hoods that are photo-consistent with others that are not.

422 Finally, the optimal depth is determined over all planes as

$$\widehat{\text{Photo}}(\mathbf{x}_j^c) = \arg \min_{Z \in [Z_{min}, Z_{max}]} \widehat{\text{Photo}}_{\pi_Z^c}(\mathbf{x}_j^c). \quad (13)$$

4. DEPTH MAP AND ALL-IN-FOCUS RECONSTRUCTION FOR FLEXIBLE SENSING 423 424

425 Algorithm 1 is presented as an example of the application of the
 426 image reprojection and photo-consistency criterion on the re-
 427 constructed 3D planes proposed in the previous section. In this
 428 algorithm, the depth maps and *all-in-focus* images on a certain
 429 3D plane can be estimated. It also shows how a reprojection is
 430 performed with two depths Z_a and Z_b . In the case of Z_a , the
 431 range of values changes as the algorithm steps forward in depth
 432 over the scene with respect to the *central camera*, assigning for

each pixel x_j^c the level of depth corresponding to the lowest photo-consistency value $\widehat{\text{Photo}}_{\pi_{Z_a}}(\mathbf{x}_j^c)$ assessed so far. In the case of Z_b , it acts like classical photo-consistency, generating planes at different levels of depth.

3 Each of the 3D points generated by reprojections \mathbf{X}_{j,Z_a}^c and \mathbf{X}_{j,Z_b}^c in the two proposed depth types is projected for each of the cameras i in the pixels \mathbf{x}_{j,Z_a}^i and \mathbf{x}_{j,Z_b}^i storing the intensities $I^i(\mathbf{x}_{j,Z_a}^i)$ and $I^i(\mathbf{x}_{j,Z_b}^i)$, respectively. Furthermore, the visibility of the projected pixels is assessed.

The main reason for the use of two levels is given by the occlusion algorithm of Wei and Quan [32], which suggests that if there is another pixel p_j^i with depth Z_a that projects to the same point in camera i as pixel \mathbf{x}_j^i and for which the projected depth is less than that of \mathbf{x}_j^i projected in that step with depth Z_b , then the pixel can be occluded. For multiple pixels warped into the same location, only the one with the smallest depth is visible, and it occludes all other projections.

Once the internal loop has finished, we can estimate the arithmetic mean of the intensities associated to the two types of depth estimation, $\bar{I}_{\pi_{Z_a}^c}$ and $\bar{I}_{\pi_{Z_b}^c}$, and their corresponding variances between the image intensities $V_{\pi_{Z_a}^c}^2$ and $V_{\pi_{Z_b}^c}^2$. With this information, we can assess (for each pixel belonging to the *central camera*) the proposed photo-consistency criterion and make a comparison, updating $\widehat{\text{Photo}}_{\pi_{Z_a}}(\mathbf{x}_j^c)$ and its depth Z_a if the photo-consistency criterion obtains a smaller value.

The algorithm satisfies the three specifications that were considered in the Introduction section: (a) establish depth planes where the objects that are at a specific scene depth are *in focus*. This is obtained by means of the arithmetic mean of the images $\bar{I}_{\pi_{Z_b}^c}$. (b) We can create an *all-in-focus* image to form the stack of depth images. This is a final by-product of the algorithm obtained when we have the final depths of the scene and project them over all the cameras. Observe that in Z_a we have stored the depth values with the lowest photo-consistency value reached until that depth plane. At the end of the loop in

$Z_b = Z_{\max}$, we have the depth of scene in Z_a , and we can obtain the arithmetic mean of the images $\bar{I}_{\pi_{Z_a}^c}$. (c) A depth map of the scene can be built with the depth stored in Z_a when the algorithm finishes.

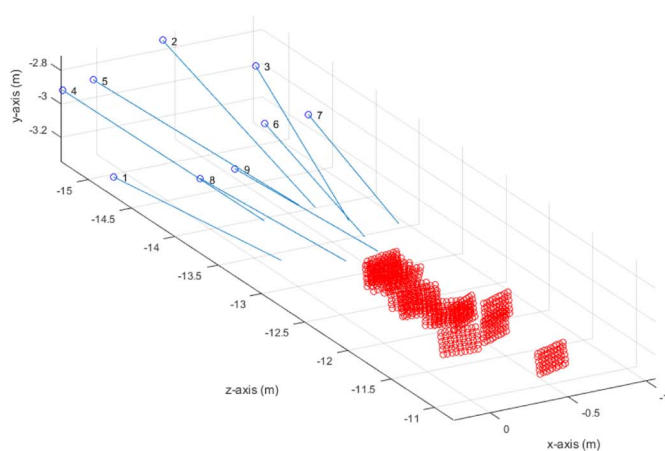
5. RESULTS

In this section, we will show some results obtained from using a flexible sensing setup. To show the capabilities of the proposed method, an experimental arrangement of the cameras, the *in focus* images obtained at different depths, and the depth map of the scene for some examples are described in the following sections.

A. Experimental Setup

Regarding the image acquisition setup, a Norpix camera array consisting of 9 AVT Mako G-192C PoE CMOS cameras (1/1.8") was used. Camera resolution was 1600 × 1200 pixels. The focal length of the optics used in the experiment was 12.5 mm. The lenses were Ricoh 12.5–75 mm F1.8, manual focus.iris/zoom lens, C-mount, 2/3" format, w/lock screws. The diagonal FOV was 39.3°. The software used for synchronized capturing was StreamPix6, for multiple camera use. The computer used to manage the entire system had a CPU Intel(R) Core(TM) i7 - 6700 K CPU at 4.0 GHz, and a speed of 2.5 GHz.

Figure 3(b) shows a picture of the experimental setup, including the array of cameras (nine cameras) and the computer used to control them. The spatial arrangement of the different cameras was located at different heights and depths, which produces a variation in the location of the objects and their size as seen by the different cameras. In addition, the cameras were positioned with arbitrary rotation to observe the scene from different points of view. As an example of what the cameras observe in the scene, we show four images (see Fig. 4). Camera number 6 [Fig. 4(c)] acted as the *central camera*, and the depth reconstruction of the scene was performed with respect to this camera. It can be observed that when choosing arbitrary posi-

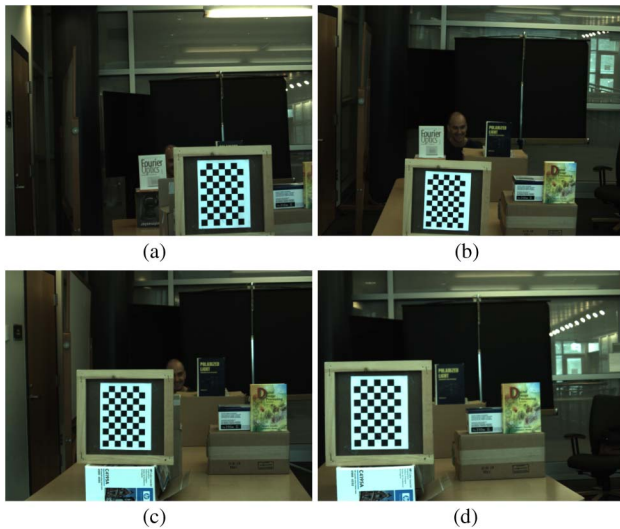


(a)



(b)

Fig. 3. (a) Centers and optical axes for the nine cameras of the setup are illustrated by solid lines in blue. Points in red represent the estimation of 3D points that belong to the movable checkerboard calibration pattern. (b) Image of camera array setup used in the experiments.



F4:1 **Fig. 4.** Elemental images: (a) camera 1, (b) camera 2, (c) camera 6,
F4:2 and (d) camera 9.

503 tions of objects at different depths in a complex scene, some
504 objects may be seen by some of the cameras and not by others.
505 This camera setup enlarges the common FOV observed by the
506 set of cameras, but it makes depth estimation of the objects in
507 the scene a difficult task.

508 B. Camera Calibration Process

509 To solve the calibration problem and therefore estimate the
510 matrix W_s , it is necessary to have a set of image points.
511 Nevertheless, it is possible that the matrix may contain some
512 missing points. Thus, the more complete the matrix is, the
513 more accurate and stable the results achieved from the calibra-
514 tion process will be. Providing a correspondence between a set
515 of scene points that can be directly visible for all cameras is a
516 difficult task. Problems arise because the inaccuracies that are
517 generated in the correspondence between the points seen by the
518 different cameras may affect the calibration accuracy. Only
519 **6** when the objects that appear in the scene have a high level
520 of texture detail is a robust correspondence likely to occur.
521 An example of this type of technique is the use of the scale-
522 invariant feature transform (SIFT), in combination with epipo-
523 lar geometry and maximum likelihood estimation in the pres-
524 ence of outliers [40].

525 In other scenes where these objects do not appear, it is pos-
526 sible to apply a moving calibration pattern in order to obtain
527 accurate information about this correspondence [41]. However,
528 the moving calibration pattern poses the same problem as the
529 direct acquisition of scene points in situations where the cam-
530 eras are far away, i.e., the correspondence points cannot be vis-
531 ible in all the views, and the partially calibrated structures have
532 to be chained together; this procedure is highly prone to errors.
533 In our case, we chose a movable checkerboard calibration pat-
534 tern, taking into account that the calibration method is able to
535 solve the existence of points that are not observed by all the
536 cameras. In Fig. 3(a), we show the camera centers and optical
537 axes for all cameras estimated during the calibration process.

Points in red represent the estimation of 3D points that belong
to the movable checkerboard calibration pattern.

7 C. Focusing Images for Different Depths

Computational reconstruction techniques in integral imaging
allow calculation of the image of the scene in a certain plane
so that the objects that are at that depth are *in focus*. Therefore,
to demonstrate the application of the proposed 3D image plane
reprojection and photo-consistency measure calculation for
flexible sensing, let us show some examples of depth maps
and the *all-in-focus* images obtained from real scenes.

In Figs. 5(a)–5(e), we show five *in-focus* images of the scene
estimated from different depth planes. To validate these im-
ages, an estimation of the distance of the object from the cam-
era array was obtained using a Laser Distance Measure Model
40-6001 to measure the depth in meters of a set of objects in
the scene [see Fig. 5(f)]. When observing the different images,
it can be seen how the object *in focus* corresponds to a part of
the scene with a sharp image while the rest of the scene is
blurred. This is because when the plane is at the depth corre-
sponding to that object, the cameras that see that object have
the same distribution of intensities, and the object is photo-
consistent at that depth.

In these demanding real experimental conditions, there is
no ground-truth available. Besides, the depth values that are
given in Fig. 5(f) are depth values taken with a laser measure
system that points only to a part of the object, from a position
close to the reference camera. Therefore, we cannot assign a
depth to a complete object, and then these depth measures
cannot be considered as measures for the object as a whole.
Nevertheless, we confirmed that the reconstruction where the
objects were in focus was the depth given by the laser measure
system.

The application of a flexible sensing setup with cameras at
arbitrary positions and different relative rotations, and therefore
with different directions for their optical axes, has both positive
and negative effects. On the one hand, the positive effect is that
by amplifying the common FOV observed by the set of cam-
eras, the size of the 3D plane that all the cameras observe can be
expanded in a larger region of the scene forming a common
mosaic. However, in this work, given that we have used a small
number of cameras for the experiments, we have chosen a
conservative configuration and limited the number of repro-
jected points in the reconstructed image planes to the number
of pixels the *central camera* has.

On the other hand, the downside effect is related to an effect
that appears in the 3D reconstruction process and that consists
in the tendency of the objects that are close to the cameras, once
they are *in focus*, to expand their corresponding boundaries in
the scene for images *in focus* at higher depths. This expansion
effect has been observed in experiments performed in classical
integral imaging from an array of cameras with parallel optical
axes and varies depending on the value of the FOV and the
distance between the cameras. This effect can be amplified
when a flexible arrangement of cameras is used, since the varia-
tion of where a close object is located by the different cameras is
bigger. As an example, the position of the checkerboard shown
in Figs. 5(c)–5(e) occupies an increasingly larger area, thus

538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594

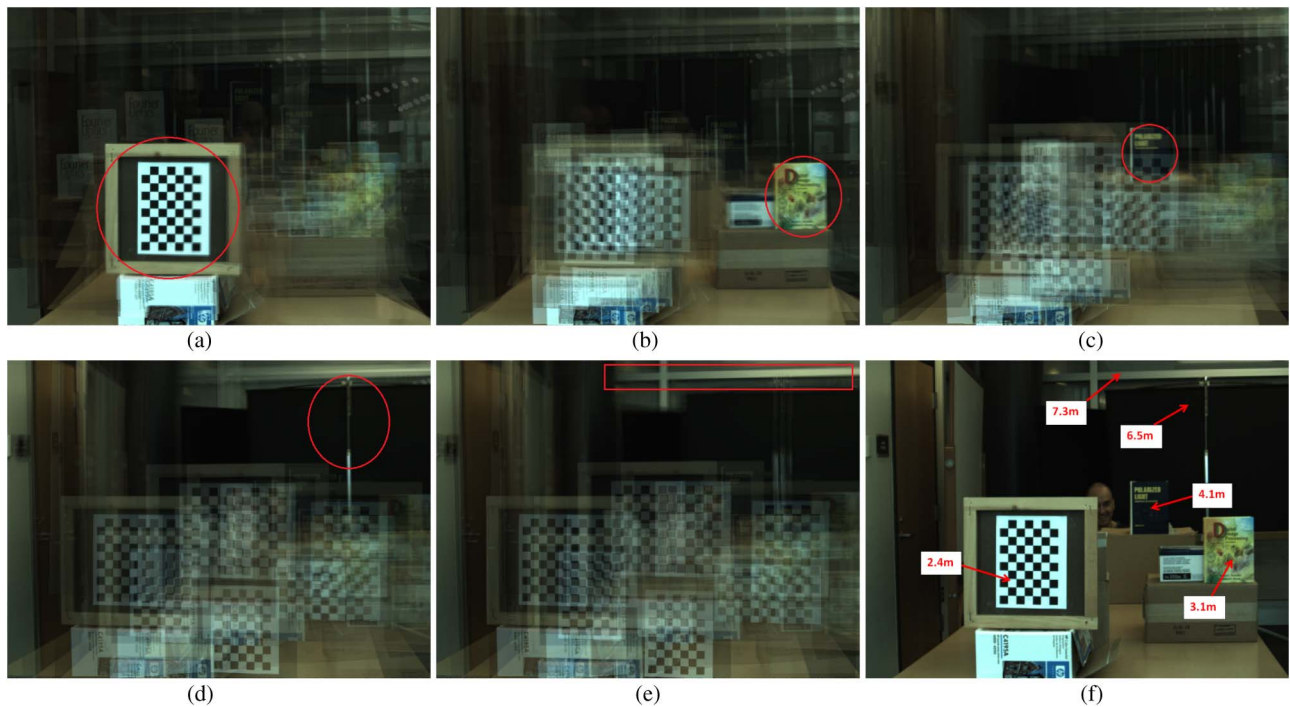


Fig. 5. (a)–(e) Reconstruction at the depth where the following objects are in focus: (a) the checkerboard; (b) first book; (c) second book; (d) the rear side of a projection screen; (e) the wall at the end of the room; (f) measured depths with the laser rangefinder.

595 interfering with objects that, for greater depths, should be seen
596 sharply.

597 **D. Depth Map and All-in-Focus Image**

598 The two other aspects addressed in this work to show how
599 integral imaging techniques can be extended to the flexible
600 sensing approach are: the generation of an *all-in-focus* image
601 and its corresponding depth map.

602 To analyze the visual quality of the proposed photo-consistency
603 criterion, the results have been compared with those
604 **8** obtained by the *Min-Var* method [21]. To do so, three scenes
605 have been analyzed, where the first one corresponds to the pre-
606 vious example of the scene that shows different objects on a
607 table, and the other two scenes basically consist of a person
608 making a gesture. Thus, two people are shown, focusing par-
609 ticularly on the reconstruction of the hand gesture and body of
610 the two subjects.

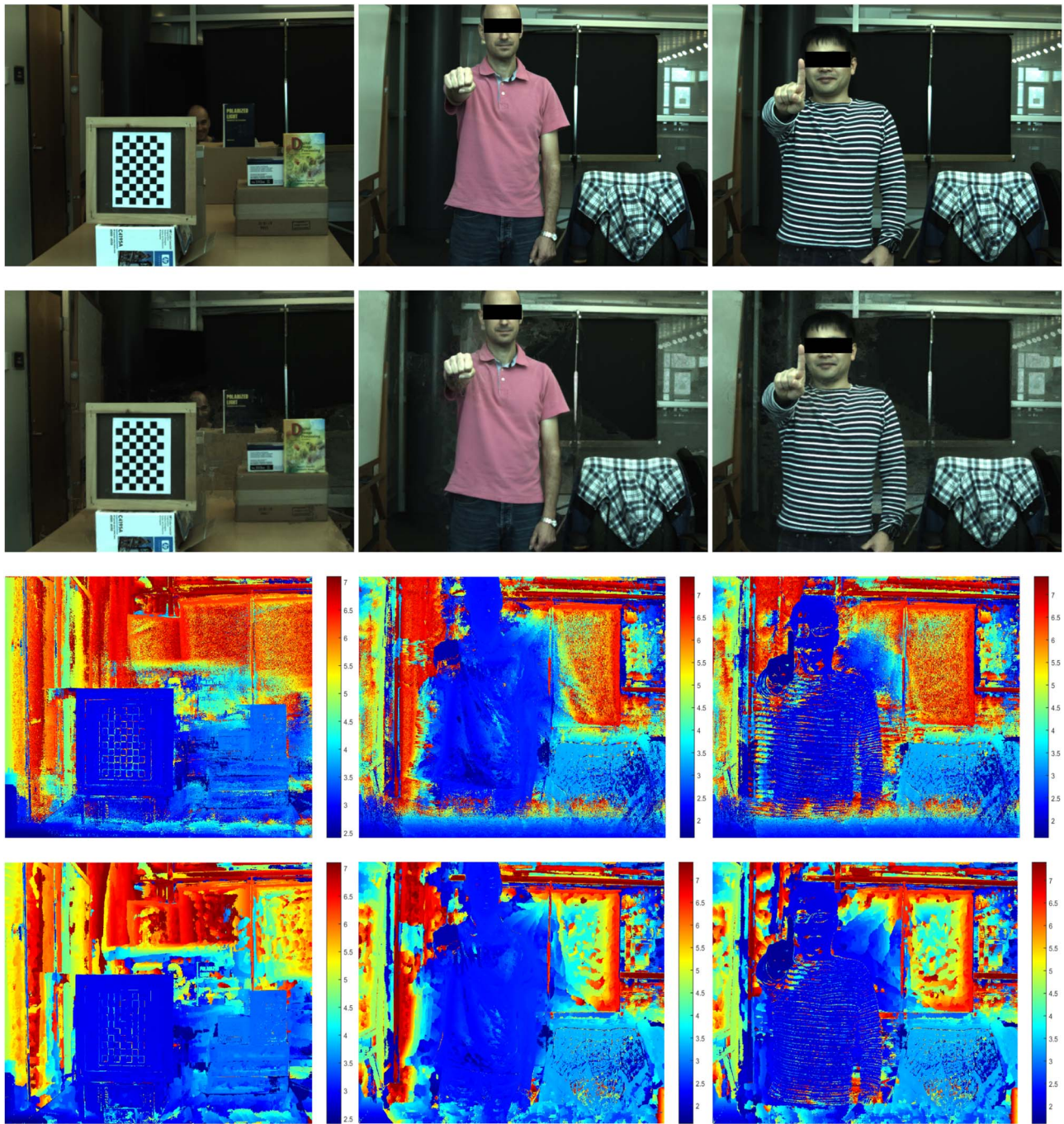
611 Figure 6 shows the results of these three scenes. The first row
612 shows the elemental images for the *central camera*. The second
613 row shows the results of the *all-in-focus* images as a function of
614 the depth estimation methodology proposed in this work. The
615 third row shows the depth map results obtained by the *Min-Var*
616 method, and the fourth row shows the depth map obtained by
617 the photo-consistency criterion used in this work.

618 The generation of an *all-in-focus* image has a strong depend-
619 ence on the photo-consistency criterion used because the im-
620 precisions that are generated during the reconstruction process
621 affect the degree of accuracy reached in the depth map. In the
622 same way, because the calculation of *all-in-focus* images is a
623 by-product of the depth map calculation, depth errors appear
624 as artifacts in the *all-in-focus* images. When comparing the

625 elemental images (first row) with their corresponding *all-in-*
626 *focus* images (second row), we may observe how some artifacts
627 and noise appear in certain regions of the scene that do not
628 allow the reconstruction to be clearly visualized. This noise
629 occurs mainly in regions that are further away from the camera
630 setup and close to objects that were in focus at a certain depth
631 and have become defocused at greater depths.

632 When we analyze the visual accuracy of the depth maps ob-
633 tained by the *Min-Var* method (see depth maps in row three of
634 Fig. 6), we must take into account that this method is based on
635 a pixel-by-pixel variance of the RGB values obtained for the
636 elemental images of each camera. We can see how this method
637 is influenced by two circumstances: the first is given by the
638 variance of intensities (called the “correspondence term” in this
639 paper) between the cameras, and it is strongly influenced by the
640 expansion effect we have previously mentioned. In addition,
641 the occlusion between objects may make it difficult to find
642 a precise correspondence of what the different cameras can ob-
643 serve. The second circumstance is given by the pixel-by-pixel
644 measurement that impedes the analysis of neighboring pixels.
645 This produces a very noisy depth map on the object surfaces.
646 Furthermore, the photo-consistency measure obtained is more
647 sensitive in the case of objects containing textured surfaces that
648 generate visual irregularity in the objects.

649 In our work, we have also added a defocus term in order to
650 measure the optimum contrast of one region of the scene that
651 is focused at a specific depth. As with the correspondence
652 term, occlusions or differences in pixel color distribution are
653 related to the point of view of the scene of each camera.
654 This fact may mean that object focusing can only be partially
655 obtained, thereby degrading the performance of this measure.



F6:1 **Fig. 6.** All-in-focus images and depth maps results. From top to bottom rows, the elemental images of the central camera, results of the all-in-focus
 F6:2 images, depth map results obtained by the Min-Var method, and depth map results obtained by the photo-consistency measure used in this work.

656 An example of this type of situation is Fig. 5(c), where objects
 657 near the checkerboard are severely affected (look at the book
 658 inside the circle). Another problem that this measure presents
 659 is the accuracy of the depth because the same object can be in
 660 focus in an interval range of depths, especially if the object has
 661 little texture.

662 When analyzing the visual results of the depth maps ob-
 663 tained with the photo-consistency measure used in this work
 664 (see depth maps in row four of Fig. 6), we can observe that
 665 the use of a bifiltering strategy based on a mean spatial filter

and a Gaussian kernel function for the intensity differences al-
 lows us to obtain results with a more homogeneous estimation
 of the object surface depth. In this case, the spatial kernel is
 centered at each individual pixel around a window W . This al-
 lows the cameras to be matched, while also adding neighborhood
 imaging information. Notice from the results that when applying
 the bifiltering process, the depth map contains lower noise, with
 smoother depth areas as a final result.

In general terms, our method is more stable in fixing
 the correct depth of the objects since it takes into account

666
 667
 668
 669
 670
 671
 672
 673
 674
 675

676 information of the neighboring pixels. Nevertheless, if depth
677 estimation is not correct, it not only affects a particular pixel,
678 but it also affects the pixels in its neighborhood. For instance,
679 we can see in rows 3 and 4 in Fig. 6 that the depth estimation
680 for the checkerboard and the books is more robust in our case
681 than for the *Min-Var* method. However, the depth estimation
682 in the black projection wall is worse in our case.

683 6. CONCLUSIONS

684 The present work has proposed a 3D reconstruction approach
685 based on integral imaging for a flexible sensing configuration of
686 the cameras. It considers that the scene is observed from notice-
687 ably different points of view in such a way that the regions of
688 the scene perceived by the cameras are difficult to match. The
689 method is based on the reprojection into 3D planes at different
690 depths that are orthogonal to the optical axis of a reference
691 camera called the *central camera*.

692 To carry out the reconstruction of the scene, a photo-
693 consistency measure combining a *defocus* and *correspondence*
694 measure has been proposed. In addition, to add information
695 from neighboring pixels, a bifiltering approach based on a mean
696 spatial filter and a Gaussian kernel function for the intensity
697 differences is applied. Based on the applied 3D plane recon-
698 struction and photo-consistency criterion, it has been shown
699 how some properties from integral imaging can be adapted
700 to this scenario. In particular, a depth estimation and an *all-*
701 *in-focus* image estimation algorithm are described to show how
702 they can be performed in a free sensing setup. Experimental
703 results show the feasibility of the proposed method and the level
704 of accuracy obtained despite the fact that the errors produced
705 by occlusions worsen in a free sensing setup. To tackle this
706 problem, an accurate multi-camera calibration method and
707 the 3D image plane reprojection approach are essential to ob-
708 tain satisfactory results.

709 The results obtained are generally consistent in real scenes
710 with different types of surfaces, although objects with a smooth
711 texture or changes due to brightness can affect the result.
712 A downside effect for objects close to the cameras is that, once
713 they are in focus, they tend to expand when reconstructing
714 planes through the scene at larger depths. This effect is ampli-
715 fied in the case of a flexible sensing configuration where the
716 optical axes are not parallel. Future work will be aimed at im-
717 proving the precision and visual quality of the generated depth
718 map, and also at incorporating other aspects such as depth map
719 regularization strategies, in an attempt to obtain smoother
720 depth maps inside homogeneous surface objects and sharp
721 estimations at depth discontinuities.

722 **9 Funding.** Generalitat Valenciana (PROMETEO-II/2014/
723 062); Spanish Ministry of Economy and Competitiveness
724 (MINECO) (ESP2013-48458-C4-3-P, MTM2013-48371-
725 C2-2-P); University Jaume I (UJIP11B2014-09); National
726 Science Foundation (NSF) (NSF/IIS-1422179); Office of
727 Naval Research (ONR) (N00014-17-1-2561).










728 **Acknowledgment.** B. Javidi acknowledges support from
729 NSF and ONR.

REFERENCES

1. S. Benton and M. Bove, *Holographic Imaging* (Wiley Interscience, 2008). 731
2. A. P. Sokolov, "Autostereoscopy and integral photography by profes- 732
sor Lippmann's method," in *Izd. MGU* (Moscow State University, 1911). 733
3. D. Scharstein and R. Szeliski, "High-accuracy stereo depth maps 734
using structured light," in *IEEE Computer Society Conference on* 735
Computer Vision and Pattern Recognition (2003), Vol. 1, pp. 195–202. 736
4. M. Hebert and E. Krotkov, "3d measurements from imaging laser ra- 737
dars: how good are they?" *Image Vis. Comput.* **10**, 170–178 (1992). 738
5. A. Bhandari and R. Raskar, "Signal processing for time-of-flight imag- 739
ing sensors: an introduction to inverse problems in computational 3-d 740
imaging," *IEEE Signal Process. Mag.* **33**(5), 45–58 (2016). 741
6. M. Martínez-Corral, A. Dorado, J. C. Barreiro, G. Saavedra, and B. 742
Javidi, "Recent advances in the capture and display of macroscopic 743
and microscopic 3-d scenes by integral imaging," *Proc. IEEE* **105**, 744
825–836 (2017). 745
7. F. Okano, J. Akai, K. Mitani, and M. Okui, "Real-time integral imaging 746
based on extremely high resolution video system," *Proc. IEEE* **94**, 747
490–501 (2006). 748
8. J. H. Park and K. M. Jeong, "Frequency domain depth filtering of 749
integral imaging," *Opt. Express* **19**, 18729–18741 (2011). 750
9. M. Daneshpanah, B. Javidi, and E. A. Watson, "Three dimensional 751
object recognition with photon counting imagery in the presence of 752
noise," *Opt. Express* **18**, 26450–26460 (2010). 753
10. D. Aloni, A. Stern, and B. Javidi, "Three-dimensional photon counting 754
integral imaging reconstruction using penalized maximum likelihood 755
expectation maximization," *Opt. Express* **19**, 19681–19687 (2011). 756
11. I. Moon and B. Javidi, "Three-dimensional visualization of objects in 757
scattering medium by use of computational integral imaging," *Opt.* 758
Express **16**, 13080–13089 (2008). 759
12. R. Schulein, C. M. Do, and B. Javidi, "Distortion-tolerant 3d recogni- 760
tion of underwater objects using neural networks," *J. Opt. Soc. Am. A* 761
27, 461–468 (2010). 762
13. B. Javidi, I. Moon, and S. Yeom, "Three-dimensional identification of 763
biological microorganism using integral imaging," *Opt. Express* **14**, 764
12096–12108 (2006). 765
14. A. Llavador, J. Sola-Pikabea, G. Saavedra, B. Javidi, and M. 766
Martínez-Corral, "Resolution improvements in integral microscopy 767
with Fourier plane recording," *Opt. Express* **24**, 20792–20798 (2016). 768
15. Y. Zhao, X. Xiao, M. Cho, and B. Javidi, "Tracking of multiple objects 769
in unknown background using Bayesian estimation in 3d space," 770
J. Opt. Soc. Am. A **28**, 1935–1940 (2011). 771
16. A. Stern and B. Javidi, "Three-dimensional synthetic aperture integral 772
imaging," *Proc. IEEE* **94**, 591–607 (2006). 773
17. J. S. Jang and B. Javidi, "Three-dimensional synthetic aperture inte- 774
gral imaging," *Opt. Lett.* **27**, 1144–1146 (2002). 775
18. M. Daneshpanah and B. Javidi, "Three-dimensional imaging with de- 776
tector arrays on arbitrary shaped surfaces," *Opt. Lett.* **36**, 600–602 777
(2011). 778
19. J. A. Rogers, T. Someya, and Y. Huang, "Materials and mechanics for 779
stretchable electronics," *Science* **327**, 1603–1607 (2010). 780
20. J. Wang, X. Xiao, and B. Javidi, "Three-dimensional integral imaging 781
with flexible sensing," *Opt. Lett.* **39**, 6855–6858 (2014). 782
21. M. Daneshpanah and B. Javidi, "Profilometry and optical slicing 783
by passive three-dimensional imaging," *Opt. Lett.* **34**, 1105–1107 784
(2009). 785
22. S. Seitz and C. Dyer, "Photorealistic scene reconstruction by voxel 786
coloring," *Int. J. Comput. Vis.* **35**, 151–173 (1999). 787
23. K. Kurulakos and S. Seitz, "A theory of shape by space carving," *Int. J.* 788
Comput. Vis. **38**, 199–218 (2000). 789
24. G. Slabaugh, W. Culbertson, T. Malzbender, M. Stevens, and R. 790
Schafer, "Methods for volumetric reconstruction of visual scenes," 791
Int. J. Comput. Vis. **57**, 179–199 (2004). 792
25. A. Martínez-Uso, P. Latorre-Carmona, J. M. Sotoca, F. Pla, and B. 793
Javidi, "Depth estimation in integral imaging based on a maximum vot- 794
ing strategy," *IEEE J. Display Technol.* (to be published). 795
26. M. Lhuillier and L. Quan, "Surface reconstruction by integrating 3d and 796
2d data of multiple views," in *IEEE Computer Society Conference on* 797
Computer Vision and Pattern Recognition (2003), Vol. 2, 1313–1320. 798
799

- 800 27. C. Strecha, T. Tuytelaars, and L. V. Gool, "Dense matching of multiple
801 wide-baseline views," in *IEEE Computer Society Conference on*
802 *Computer Vision and Pattern Recognition* (2003), Vol. 2, 1194–1201.
803 28. C. Strecha, R. Fransens, and L. V. Gool, "Wide-baseline stereo from
804 multiple views: a probabilistic account," in *IEEE Computer Society*
805 *Conference on Computer Vision and Pattern Recognition* (2004),
806 Vol. 2, 552–559.
807 29. Y. Furukawa and C. Hernandez, "Multi-view stereo: a tutorial," *Found.*
808 *Trends Comput. Graph. Vis.* **9**, 1–148 (2015).
809 30. T. Svoboda, D. Martinec, and T. Pajdla, "A convenient multi-camera
810 self-calibration for virtual environments," *Presence Teleop. Virt.*
811 *Environ.* **14**, 407–422 (2005).
812 31. T.-C. Wang, A. A. Efros, and R. Ramamoorthi, "Depth estimation with
813 occlusion modeling using light-field cameras," *IEEE Trans. Pattern*
814 *Anal. Mach. Intell.* **38**, 2170–2181 (2016).
815 32. Y. Wei and L. Quan, "Asymmetrical occlusion handling using graph cut
816 for multi-view stereo," in *IEEE Computer Society Conference on*
817 *Computer Vision and Pattern Recognition* (2005), Vol. 2, pp. 902–909.
818 33. S. M. Seitz, B. Curless, J. Diebel, D. Scharstein, and R. Szeliski, "A
819 comparison and evaluation of multi-view stereo reconstruction algo-
820 rithms," in *IEEE Computer Society Conference on Computer Vision*
821 *and Pattern Recognition* (2006), Vol. 1, pp. 519–526.
34. R. Harley and A. Zisserman, *Multiple View Geometry in Computer*
822 *Vision* (Cambridge University, 2000). 823
35. P. Sturm and B. Triggs, "A factorization based algorithm for multi-
824 image projective structure and motion," in *European Conference on*
825 *Computer Vision* (Springer-Verlag, 1996), pp. 709–720. 826
36. D. Martinec and T. Pajdla, "Structure from many perspective images
827 with occlusions," in *European Conference on Computer Vision*, A.
828 Heyden, G. Sparr, M. Nielsen, and P. Johansen, eds. (Springer-
829 Verlag, 2002). 830
37. G. Wang and Q. M. J. Wu, *Guide to Three Dimensional Structure and*
831 *Motion Factorization* (Springer-Verlag, 2011). 832
38. R. Szeliski and P. Golland, "Stereo matching with transparency and
833 matting," *Int. J. Comput. Vis.* **32**, 45–61 (1999). 834
39. H. Saito and T. Kanade, "Shape reconstruction in projective grid
835 space from large number of images," in *IEEE Computer Society*
836 *Conference on Computer Vision and Pattern Recognition* (1999),
837 pp. 49–54. 838
40. P. H. S. Torr and A. Zisserman, "MLESAC: a new robust estimator
839 with application to estimating image geometry," *Comput. Vis.*
840 *Image Underst.* **78**, 138–156 (2000). 841
41. Z. Zhang, "A flexible new technique for camera calibration," *IEEE*
842 *Trans. Pattern Anal. Mach. Intell.* **22**, 1330–1334 (2000). 843

Queries

1. AU: Please check my edits here: “(c) the ability to infer a depth map of the scene” the original was “(c) to infer a depth map of the scene,” which didn’t quite fit with the grammatical structure of (a) and (b). 
2. AU: Please check my edit in the sentence beginning, “Thus, we wish to estimate the projection matrices...” I changed “which” to “that.” 
3. AU: Please check my edits in the sentence beginning, “Each of the 3D points generated...” 
4. AU: Please provide value (8”) in SI unit instead of “inches.” 
5. AU: Please check my edit in the sentence beginning, “ The computer used to manage the entire system...” 
6. AU: Please check my edits in the sentence beginning, “Only when the objects that appear in...” 
7. AU: Please check my edits in the sentence beginning, “Therefore, to demonstrate the application...” 
8. AU: Does *Min-Var* have to be capped and italicized, or could it be min-var? (I didn’t change it.) 
9. AU: The funding information for this article has been generated using the information you provided to OSA at the time of article submission. Please check it carefully. If any information needs to be corrected or added, please provide the full name of the funding organization/institution as provided in the CrossRef Open Funder Registry (<http://www.crossref.org/fundingdata/registry.html>). 
10. AU: Is updated information available for Ref. [25]? 