



High-accuracy 3D locators tracking in real time using monocular vision

C. Elmo Kulanesan¹ · P. Vacher¹ · L. Charleux¹ · E. Roux¹

Received: 2 March 2023 / Revised: 29 August 2023 / Accepted: 28 November 2023 / Published online: 11 January 2024
© The Author(s), under exclusive licence to Springer-Verlag GmbH Germany, part of Springer Nature 2024

Abstract

In the field of medical applications, precise localization of medical instruments and bone structures is crucial to ensure computer-assisted surgical interventions. In orthopedic surgery, existing devices typically rely on stereoscopic vision. Their purpose is to aid the surgeon in screw fixation of prostheses or bone removal. This article addresses the challenge of localizing a rigid object consisting of randomly arranged planar markers using a single camera. This approach is especially vital in medical situations where accurate object alignment relative to a camera is necessary at distances ranging from 80 cm to 120 cm. In addition, the size limitation of a few tens of centimeters ensures that the resulting locator does not obstruct the work area. This rigid locator consists of a solid at the surface of which a set of plane markers (ArUco) are glued. These plane markers are randomly distributed over the surface in order to systematically have a minimum of two visible markers whatever the orientation of the locator. The calibration of the locator involves finding the relative positions between the individual planar elements and is based on a bundle adjustment approach. One of the main and known difficulties associated with planar markers is the problem of pose ambiguity. To solve this problem, our method lies in the formulation of an efficient initial solution for the optimization step. After the calibration step, the reached positioning uncertainties of the locator are better than two-tenth of a cubic millimeter and one-tenth of a degree, regardless of the orientation of the locator in space. To assess the proposed method, the locator is rigidly attached to a stylus of about twenty centimeters length. Thanks to this approach, the tip of this stylus seen by a 16.1 megapixel camera at a distance of about 1 m is localized in real time in a cube lower than 1 mm side. A surface registration application is proposed by using the stylus on an artificial scapula.

Keywords Monocular vision · 3D locators · Planar marker · 3D tracking · Aruco

1 Introduction

The pose estimation of an object in a scene is a classical computer vision problem. It is generally seen as a least-squares optimization problem called Perspective-N-Points (PNP) in which the reprojection errors of a number of points of interests (POIs) located on a target and observed on several images are minimized [1–3]. For real-time applications, the use of easy-to-detect fiducial markers provides excellent computing time performance to detect them on the images. At the same time, the use of a single camera, also referred to as monocular vision, reduces financial costs as well as the complexity of the setup. Many applications presented in the literature take place in this context ranging from positioning for augmented reality, drone navigation, gesture recognition[4–7].

The fiducial markers used in the literature are mostly planar and use the four corners of a square as POIs. While four coplanar POIs are theoretically sufficient for a pose estimation based on a single marker, in practice, the PNP problem becomes ill-posed. Two different solutions of pose estimation may exist, especially at large working distance and large inclination. Figure 1 presents two configurations illustrating this problem of pose ambiguity [8] in a simple case, rotation around the y axis of the marker. The first one (*configuration A*) corresponds to a non-ambiguous observation of the pose, the observed projections allow without difficulty to establish a reliable pose using the classical pose estimation algorithms. In the second configuration, the marker is far from the camera. We observe that two different orientations of the marker lead to an almost identical projection on the optical sensor. In our application, the test conditions correspond to this second configuration. In the absence of measurement noise on the detection of the four corners of the marker, the two poses are unambiguously discernible. In practice, measurement noise

✉ C. Elmo Kulanesan
christian.elmo-kulanesan@univ-smb.fr

¹ Université Savoie Mont Blanc Laboratoire SYMME, 74940
Annecy-le-Vieux, France

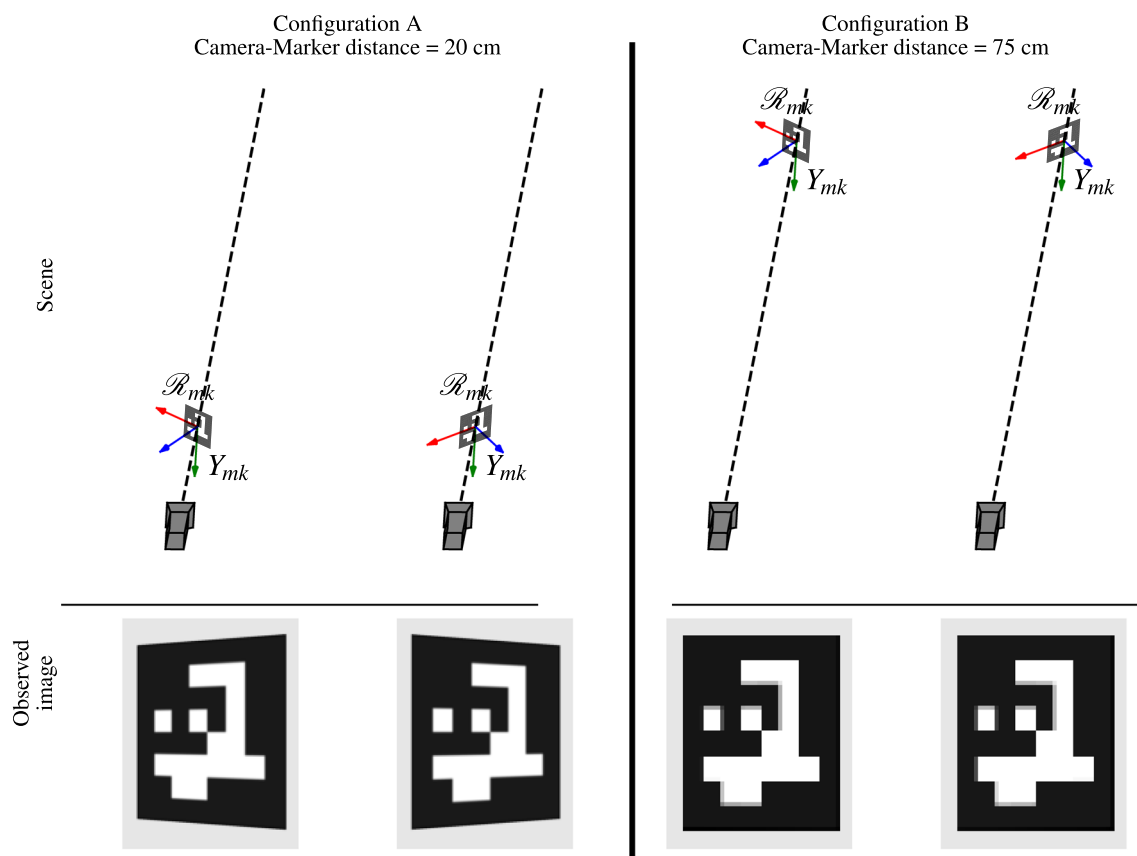


Fig. 1 Illustration of the pose ambiguity phenomenon. Two configurations are presented where a marker of 12 mm side is seen by a camera whose sensor is located at two distinct distances. In this illustration, the camera has a resolution of 2048 px \times 1536 px and a sensor size of 7.07 mm \times 5.3 mm. The focal length is 12 mm. In each configuration, two orientations of the marker differing from a simple rotation θ around its axis Y_m by $\pm 36^\circ$ are presented. In configuration A, the markers are positioned at a distance of 20 cm from the camera. The projections on

the images clearly show the perspective effect of the markers. Visually, the vertical edges are of different lengths. In configuration B, the markers are positioned at a distance of 75 cm from the camera. The images observed for these two situations are almost identical despite the very different orientations of the two markers. A small uncertainty on the identification of the spatial position of the 4 points of this quadrilateral can lead to strongly different pose estimation of this marker at the origin of the pose ambiguity

exists and the solution obtained by optimization does not always correspond to the real pose of the marker. The global approach proposed here consists in making the placement of the locator more reliable by increasing the number of markers visible on the images.

1.1 Related work

The related work in the field of pose estimation and localization spans various domains, including monocular and multiview camera systems as well as algorithms both with and without dedicated markers. In the realm of single-view pose estimation, the DTAM framework proposed by Newcombe et al. [9] enables real-time dense tracking and mapping, while Kendall and Cipolla [10] address uncertainty modeling in deep learning for camera relocalization. Expanding to multiview scenarios, the seminal work by Hartley and

Zisserman [11] outlines the fundamentals of multiple view geometry, and Snavely et al.'s concept of "photo tourism" [12] introduces the notion of 3D exploration through photo collections.

In the absence of dedicated markers, markerless approaches have gained prominence. Newcombe et al.'s KinectFusion [13] facilitates real-time dense surface mapping and tracking through a depth sensor, and Mur-Artal and Tardós's ORB-SLAM2 [14] presents a versatile open-source SLAM system catering to various camera setups.

Notably, optimization plays a pivotal role in refining these methods. The Levenberg–Marquardt algorithm [15] has been a cornerstone, as Levenberg's method [16] and Marquardt's algorithm [17] significantly contribute to solving nonlinear least-squares problems. Additionally, bundle adjustment, as elaborated by Triggs et al. [18], and alternative optimization techniques such as those detailed by Lepetit and Lourakis

[19, 20] continue to enhance the precision and efficiency of pose and localization estimation algorithms by working on a better mathematical representation of Rodrigues' parameters.

Some studies proposed solutions to solve the pose ambiguity by using a fifth or more POIs [21, 22]. In some configurations, these approaches show low reliabilities, as demonstrated by [23].

In a noise-free configuration, this second solution exhibits a higher reprojection error than the true solution and can therefore be easily discarded [8]. In practice, it is common for the bad solution to have a lower reprojection error than the good one and makes it impossible to differentiate them directly [24]. Consequently, a single plane marker cannot always be used to estimate the pose of an object with confidence.

To avoid pose ambiguity, several markers can be positioned in a non-coplanar way on the object [25, 26]. In this paper, this set of markers stuck on a rigid object is called *locator*. A calibration is needed if the relative position of the markers is not precisely controlled. The calibration of a locator can be seen as a Structure-from-Motion (SfM) problem [12, 27] which can be solved offline to obtain an accurate model of the relative position of the fiducial markers within the locator [28, 29].

As before, this is an ill-posed problem which requires an accurate enough starting point in order to converge toward the solution. The pose ambiguity can jeopardize the establishment of a well-suited starting point. Thus, in some applications, the markers must be placed according to a predefined pattern (e.g., on a plane grid or on the faces of a polyhedron) in order to make a direct resolution of the PNP problem possible [26, 30, 31]. In the more general case considered here where the arrangement of markers is random, a specific method must be developed to build an initial solution taking into account the pose ambiguity.

1.2 Our contributions

In this paper, we focus on the problem of ambiguity of the pose as well as the identification of markers that are incorrectly detected or damaged but still detected (e.g., a bent marker or partially occluded marker) as shown in Fig. 2. We propose a new method to construct a reliable estimation of the relative position of fiducial markers. This method allows the elimination of ambiguous positions as well as bad detections. The paper is organized as follows. Section 2 presents the developed method after laying the mathematical foundations. Section 3 describes a series of experiments setup to demonstrate the performance of the proposed method. Section 5 proposes a discussion and a conclusion regarding the results and the overall performance of the method.

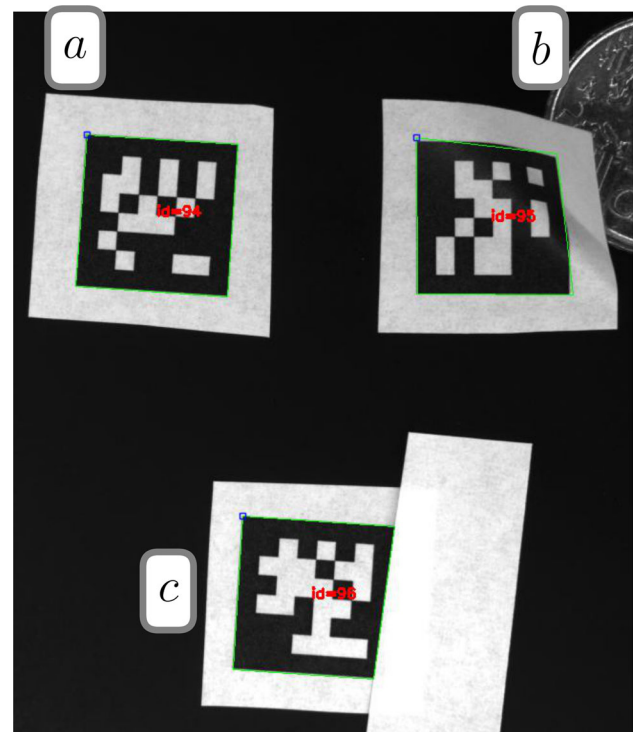


Fig. 2 Illustration of two configurations that can lead to incorrect pose estimation. **a** Appropriate detection of the marker. **b** Bad marker: mis-detection of the marker. In this case, it is bent but a stain could produce the same effect. The associated poses are systematically wrong. **c** Bad detection: a marker is detected incorrectly. In this case, the marker is partially occluded but is still detected. One of its corners is detected in the wrong position and the associated poses are erroneous

2 Proposed method

2.1 Initial concepts and definitions

2.1.1 Definitions

In our approach, a set of planar square fiducial markers M_i with $i \in \{1, \dots, n_M\}$ are randomly glued at the surface of an object with an arbitrary geometry. It is assumed that the markers are not in a coplanar position. This object, noted L , is referred to as the locator.

In the present work, ArUco markers [32] are used, but this approach could be similar to all type of square planar fiducial markers such as AprilTag [33] or ARTag [34].

A set of images I_j of the locator with $j \in \{1, \dots, n_I\}$ is created from different viewpoints with a single camera. Each marker M_i and each image I_j are associated with a reference frame. The images are processed with the ArUco method in order to detect the visible markers [35, 36]. The pixel coordinates of the 4 corners of a marker M_i visible on the image I_j are noted $\tilde{u}_{M_i/\mathcal{I}_j}(I_j) \in \mathbb{R}^2 \times \mathbb{R}^4$. The sub-

pixel position of these corners is refined using the ‘‘AprilTag 2 method’’ [37].

2.1.2 Affine transformations

When passing from one reference frame to another, the coordinates of a point undergo a rotation \mathbf{r} followed by a translation \mathbf{t} . For the sake of simplicity, we use notations consistent with the ones used by [24]. Accordingly, such an affine transformation is noted $\gamma = (\mathbf{r}, \mathbf{t})$, with $\mathbf{r} = [r_x, r_y, r_z]^T$, the rotation vector as defined by Rodrigues convention and $\mathbf{t} = [t_x, t_y, t_z]^T$, the translation vector. The rotation matrix \mathbf{R} , bijectively associated to \mathbf{r} verifies:

$$\mathbf{R} = \mathbf{I} + \bar{\mathbf{r}} \sin \theta + \bar{\mathbf{r}}^2(1 - \cos \theta) \tag{1}$$

Where $\theta = \|\mathbf{r}\|$ and $\bar{\mathbf{r}}$ is the cross-product matrix:

$$\bar{\mathbf{r}} = \frac{1}{\theta} \begin{bmatrix} 0 & -r_z & r_y \\ r_z & 0 & -r_x \\ -r_y & r_x & 0 \end{bmatrix}, \tag{2}$$

and $\theta \in [0, 2\pi]$, due to the periodicity properties of rotations. The affine transformation matrix $\Gamma(\gamma) \in \mathbb{R}^4 \times \mathbb{R}^4$ is then defined by:

$$\Gamma(\gamma) = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ 0 & 0 & 0 & 1 \end{bmatrix} \tag{3}$$

Affine transformation matrices $\Gamma_a = \Gamma(\gamma_a)$ and $\Gamma_b = \Gamma(\gamma_b)$ can be composed by matrix product so that $\Gamma_{b,a} = \Gamma_b \otimes \Gamma_a$. By extension, it is assumed that the symbol (\cdot) indicates the composition for γ transformations so that $\gamma_{b,a} = \gamma_b \cdot \gamma_a$.

2.1.3 Single marker pose estimation

In its own reference frame \mathcal{R}_{M_i} , the spatial homogeneous coordinates of the corners s_k with $k \in \{0, \dots, 3\}$ of a marker M_i are consolidated in matrix $\mathbf{C}_{M_i/\mathcal{R}_{M_i}}$ such as:

$$\mathbf{C}_{M_i/\mathcal{R}_{M_i}} = \frac{1}{2} \begin{matrix} s_0 & s_1 & s_2 & s_3 \\ \uparrow & \uparrow & \uparrow & \uparrow \\ \begin{bmatrix} -d & d & d & -d \\ -d & -d & d & d \\ 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix} \end{matrix}, \tag{4}$$

with d the side length of each marker. These coordinates can be expressed in any reference frame \mathcal{R} by their homogeneous coordinates:

$$\mathbf{C}_{M_i/\mathcal{R}} = \Gamma_{\mathcal{R} \leftarrow \mathcal{R}_{M_i}} \otimes \mathbf{C}_{M_i/\mathcal{R}_{M_i}} = \gamma_{\mathcal{R} \leftarrow \mathcal{R}_{M_i}} \cdot \mathbf{C}_{M_i/\mathcal{R}_{M_i}} \tag{5}$$

By composing the resulting coordinates of the corners of the marker M_i with a Ψ operator, these points are projected into the pixel space of the camera on the I_j image. Ψ is a unique affine transformation obtained by calibrating the camera sensor which integrates the change of reference frame from the camera frame \mathcal{R}_c to the image frame \mathcal{R}_I by perspective projection and a change of scale operation. The coordinates of the projected points will then be expressed as:

$$\mathbf{u}_{M_i/\mathcal{R}_I}(\gamma(I_j)) = \Psi(\delta) \cdot \gamma(I_j) \cdot \mathbf{C}_{M_i/\mathcal{R}_{M_i}} \tag{6}$$

Where $\delta = (f_x, f_y, c_x, c_y, k_1, \dots)$ is the tuple of intrinsic camera parameters, (f_x, f_y) are the focal lengths parameters, c_x, c_y the optical center parameters and (k_1, \dots) the distortion parameters. The reprojection error of the marker M_i on the image I_j , as a function of $\gamma(I_j)$, is then written:

$$e_{M_i, I_j}(\gamma(I_j)) = \sum_{m=1}^4 \sum_{l=1}^2 \left[\Delta \mathbf{u}_{M_i/\mathcal{R}_I}(\gamma(I_j)) \right]_{l,m}^2, \tag{7}$$

with l and m denote, respectively, the number of coordinates and the number of corners of the marker M_i . $\Delta \mathbf{u}_{M_i/\mathcal{R}_I} \in \mathbb{R}^2 \times \mathbb{R}^4$ are the reprojection residues in pixels calculated from the difference between the measured corners $\tilde{\mathbf{u}}_{M_i/\mathcal{R}_I}(I_j)$ and the projected corners $\mathbf{u}_{M_i/\mathcal{R}_I}(\gamma(I_j))$:

$$\Delta \mathbf{u}_{M_i/\mathcal{R}_I}(\gamma(I_j)) = \mathbf{u}_{M_i/\mathcal{R}_I}(\gamma(I_j)) - \tilde{\mathbf{u}}_{M_i/\mathcal{R}_I}(I_j) \tag{8}$$

In theory, the pose of the marker M_i relative to the camera reference frame \mathcal{R}_c by an image I_j can be determined by minimizing the reprojection error according to $\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_i}}$:

$$\begin{aligned} \gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_i}} &= \arg \min_{\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_i}} \in \mathbb{R}^3 \times \mathbb{R}^3} \\ e_{M_i, I_j}(\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_i}}) \end{aligned} \tag{9}$$

The knowledge of the pose of each individual marker allows the locator structure to be determined. This structure is retrieved by an optimization step called bundle adjustment described in the paragraphe 2.1.6. In practice, some pose estimation of individual marker can be wrong, due to pose ambiguity or misdected markers.

2.1.4 Pose ambiguity

Due to the coplanarity of the corners of each individual marker, the reprojection error $e_{M_i, I_j}(\gamma(I_j))$ can have two minimums [See [3], Fig. 1–3]. When two solutions exist, they are noted $\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_i}, 0}$ and $\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_i}, 1}$. One is the real solution and the other is a bad solution. The presence of measurement noise implies that in some cases, the

minimum of the reprojection error does not necessarily correspond to the real solution. This problem is called *pose ambiguity* [8, 38]. In the following, we systematically calculate these two poses using the method proposed by [3]. If they are different, both are kept and the right one will be identified afterward. We used markers of size $d = 12.41$ mm at a typical working distance of $L = 800$ to 1200 mm. In this configuration, the probability of presence of two poses is almost systematic.

2.1.5 Bad markers and bad detections

In practice, the following two configurations can also lead to incorrect pose estimation illustrated in Fig. 2:

Bad marker: A damaged marker but still detected. It occurs when a marker is no longer flat or if it is stained or scratched. In this case, both poses are systematically incorrect.

Bad detection: Poor marker detection typically caused by a corner being detected at the wrong position. This can happen on partially occluded markers or bad image quality. In this case, the detection results also in two incorrect poses.

The detections from these configurations must be eliminated before the bundle adjustment to allow convergence to the correct solution. This implies eliminating the bad detections and not considering the bad markers in the construction of the locator.

2.1.6 Bundle adjustment

The pose and structure of the locator L are determined by solving a bundle adjustment problem. Since the markers are not coplanar, the locator is not affected by pose ambiguity for a single marker. Firstly, among all the markers, a marker M_0 is arbitrarily taken as reference within the locator. The structure of the locator will be fully defined by knowledge of $\mathcal{Y}_S = \{\gamma_{\mathcal{R}_{M_0} \leftarrow \mathcal{R}_{M_i}} | i \in \{1, \dots, n_M\}\}$. The pose of the locator in the image batch I_j is defined by:

$$\mathcal{Y}_I = \{\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_0}} | j \in \{1, \dots, n_I\}\} \tag{10}$$

The reprojection error is defined as:

$$E(\mathcal{Y}_S^*, \mathcal{Y}_I^*) = \sum_{i=1}^{n_M} \sum_{j=1}^{n_I} \eta_{M_i, I_j}(\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_0}} \cdot \gamma_{\mathcal{R}_{M_0} \leftarrow \mathcal{R}_{M_i}}), \tag{11}$$

where:

$$\eta_{M_i, I_j}(\gamma) = \begin{cases} e_{M_i, I_j}(\gamma), & D_{M_i, I_j} = 1 \\ 0, & D_{M_i, I_j} = 0 \end{cases}, \tag{12}$$

and D a detection mask defined as follows:

$$D_{M_i, I_j} = \begin{cases} 1, & \exists(i, j) \mid M_i \subset I_j \\ 0, & \text{otherwise} \end{cases} \tag{13}$$

Consequently, the optimal parameters \mathcal{Y}_S and \mathcal{Y}_I are the solutions of:

$$\mathcal{Y}_S, \mathcal{Y}_I = \arg \min_{\mathcal{Y}_S^*, \mathcal{Y}_I^*} E(\mathcal{Y}_S^*, \mathcal{Y}_I^*) \tag{14}$$

The optimal result highly depends on the initial guess provided $(\mathcal{Y}_{S,0}, \mathcal{Y}_{I,0})$. We propose a method based on graph theory that is able to handle several problems prior to the bundle adjustment. It is therefore essential to eliminate bad poses beforehand by removing the pose ambiguity and by eliminating bad detections of markers.

2.2 Proposed algorithm

2.2.1 Detection graph and cycle basis

We consider a graph \mathcal{G}_S whose vertices are the markers and the images. The sets of images and markers being independent, each detection $D_{M_i, I_j} = 1$ corresponds to an edge of this graph. A cycle basis $\{\xi_i | i \in \{1, \dots, N_c\}\}$ is then generated from this graph [39]. A basic example is shown in Fig. 3.

2.2.2 Pose classification using cyclic rotational errors

Let us consider the cycle $\xi_0 = (M_1, I_1, M_2, I_0)$ from Fig. 3 with length $N_e = 4$. For each edge of the cycle corresponding to the detection of the marker M_i on the image I_j , the true affine transformation is noted $\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_i}}$. By composing these transformations, a residual pose $\hat{\gamma}_{\xi_0} = (\hat{\mathbf{r}}_{\xi_0}, \hat{\mathbf{t}}_{\xi_0})$ for this cycle is calculated as:

$$\hat{\gamma}_{\xi_0} = \gamma(I_1)_{\mathcal{R}_{M_2} \leftarrow \mathcal{R}_c} \cdot \gamma(I_0)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_1}} \cdot \gamma(I_1)_{\mathcal{R}_{M_1} \leftarrow \mathcal{R}_c} \cdot \gamma(I_0)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_2}} \tag{15}$$

If the pose estimations were free of errors, this residual pose would be null. In practice, there is always some degree of error in the pose estimations. In addition, in case of pose ambiguity, there are two pose estimations $\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_i}}, 0$ and $\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_i}}, 1$. It follows that the residual affine

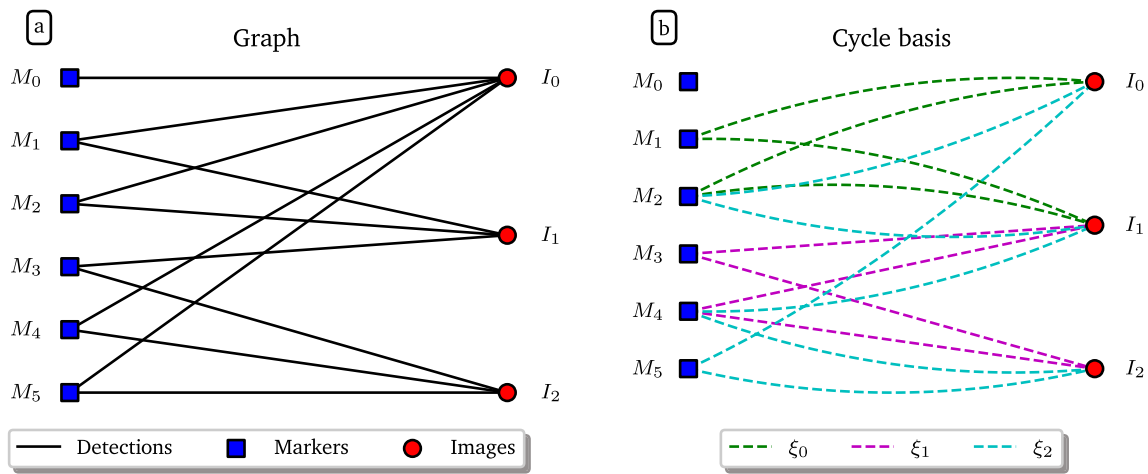


Fig. 3 Simple example of construction of a graph \mathcal{G}_s . **a** The sets of three images and six markers are represented. Each detection D_{M_i, I_j} of a marker M_i on an image I_j corresponds to an edge of the graph. For example, the markers M_1, M_2 and M_3 are here detected on the image

I_1 . **b** The cycle basis allows to browse at least once each edge. Here, three cycles ξ_i are necessary to build the base. The marker M_0 was observed only once, so it does not belong to any cycle and is therefore excluded from the cycle base

transformation $\hat{\gamma}_{\xi_0}^k$ associated with a given cycle can be computed according to N_p combinations of estimated poses for $k \in \{1, \dots, N_p\}$ where $1 \leq N_p \leq 2^{N_e}$. *A priori*, only one of these combinations should be correct while the others should produce large residual affine transformations. Since the pose ambiguity predominantly affects the rotational component of a transformation γ , the angular residual $\hat{\theta}_{\xi_0}^k = \|\hat{\gamma}_{\xi_0}^k\|$ is used as the sole indicator of combination quality. The following criterion is therefore used to validate the combinations associated with a cycle of length N_e :

$$\frac{\hat{\theta}_{\xi_0}^k}{N_e} \leq \hat{\theta}_c, \tag{16}$$

where $\hat{\theta}_c$ is the critical angle below which the angular error per edge is considered negligible. This angle is arbitrarily set up to 0.1° in our application.

In fact, we observe that several combinations of the same cycle may often validate this criterion. This implies that incorrect pose estimations can be involved in combinations that validate the criterion of Eq. 16. There are two main reasons for this observation. First, the presence of almost identical images within the same cycle necessarily implies that two combinations of the cycle will validate the angular criterion. Secondly, fortuitously and due to the large number of pose estimates, we sometimes observe validations of this angular criterion for wrong cycles. A statistical approach can separate the correct pose estimations from the incorrect ones.

A counter $C_{I_j, M_i, k}$ is associated with each pose $\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_i}, k}$ for $k \in \{0, 1\}$, which can belong to one of the following classes: bad (B), undetermined (U) and good (G). The goal of our method is to move the poses from their

initial class (U) to classes (B) and (G). All the combinations of each cycle are tested using Eq. 16 and each time, one of them is validated, the counters are incremented by 1. We assume that incorrect validations of combinations involving bad poses are in the minority. We thus assume that of two poses associated with the same detection, the good one will be validated at least α times more often than the bad one. The α criterion is formalized by the following equation:

$$\alpha = \frac{\max(C_{I_j, M_i, 0}, C_{I_j, M_i, 1})}{\min(C_{I_j, M_i, 0}, C_{I_j, M_i, 1})} \tag{17}$$

To secure the validation of good poses, the value $\alpha = 1.5$ has been chosen.

Table 1 proposes an example of implementation of this approach. The final value of this counter allows to treat each detection associated with poses belonging to (U) individually. In the absence of pose ambiguity, if $C_{I_j, M_i, k} > 0$, the only pose is reclassified as (G), (B) otherwise. If $C_{I_j, M_i, 0}/C_{I_j, M_i, 1} \geq \alpha$, the pose $\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_i}, 0}$ is considered as (G) and $\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_i}, 1}$ (B). Conversely, if $C_{I_j, M_i, 1}/C_{I_j, M_i, 0} \geq \alpha$, the pose $\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_i}, 1}$ is considered (G) and $\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{M_i}, 0}$ (B). Finally, in the remaining cases, neither pose is clearly better and both are kept in (U). The graph \mathcal{G}_s is updated from the detections present in the classes (G) and (U) and a new basis cycle is generated. The $C_{I_j, M_i, k}$ are reset to 0 for all the remaining poses. Equation 16 is once again applied to increment the counters.

This procedure is re-iterated until the size of the (U) class does not decrease during an iteration. If the graph obtained during the iterations is disjoint then the iterations are stopped.

Table 1 Example of application of the proposed method in the case of the graph of Fig. 3

Id.	Detection		Valid. counter		To class	
	Mark.	Img.	γ_0	γ_1	γ_0	γ_1
0	M_1	I_0	1	0	G	B
1	M_1	I_1	2	2	U	U
2	M_2	I_0	2	8	B	G
3	M_2	I_1	5		G	
4	M_3	I_1	0		B	
...						

Only the poses which are in class (U) at the beginning of the current iteration are represented. The column "To class" indicates to which class the poses are moved. Detections 3 and 4 have no pose ambiguity which explains the absence of γ_1 in these two cases. We can see that except for the poses associated with detection 1, all the poses are initially in (U) are placed in (B) and (G)

The remaining poses in class (U) are then placed in (B) then all poses are classified as (G) or (B).

Consequently, each edge of the final graph \mathcal{G}_s is associated to a single pose belonging to the class (G).

2.2.3 Formalization of the educated guess initial solution

The formulation of the initial solution requires the selection of a reference marker from which the locator structure is described. Among the markers, the one with the lowest eccentricity is chosen as the reference marker and noted M_0 as in Sect. 2.1.6.

The relative pose of each marker $\gamma_{S,0}$ and each frame $\gamma_{I,0}$ is then estimated by looking for the shortest path between the reference marker and each other node of the simple graph \mathcal{G}_s [40].

The main interest of this algorithm lies in the quality of this initial solution. Indeed, it is close enough to the optimal solution and thus represents an excellent starting point for a bundle adjustment procedure described in Sect. 2.1.6 which allows to further refine this solution to obtain γ_S and γ_I from Eq. 14.

2.3 Locator pose estimation from a single image

After the proposed structural calibration, the pose of the locator can be estimated from a single image. The coordinates of the marker corners are then known relatively to the chosen reference marker M_0 . The locator's reference frame \mathcal{R}_L is then equal to the reference frame \mathcal{R}_{M_0} . Knowing the locator structure γ_S , the coordinates of the corners of the markers are expressed in the locator's reference frame \mathcal{R}_L such as:

$$C_{M_i|\mathcal{R}_L} = \gamma_S \cdot C_{M_i|\mathcal{R}_{M_i}} \tag{18}$$

Consequently, the coordinates of the markers in the camera reference frame \mathcal{R}_c are given by the following expression:

$$C_{M_i|\mathcal{R}_c}(I_j) = \gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_L} \cdot C_{M_i|\mathcal{R}_L} \tag{19}$$

The coordinates of the projections of the corners of the markers visible to the camera ($D_{M_i,I_j} = 1$) projected onto the image I_j can be expressed in the image reference frame \mathcal{R}_I as:

$$u_{M_i|\mathcal{R}_I}(I_j) = \Psi(\delta) \cdot \gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_L} \cdot C_{M_i|\mathcal{R}_L} \tag{20}$$

The reprojection error residuals $\Delta u_{M_i|\mathcal{R}_I} \in \mathbb{R}^2 \times \mathbb{R}^4$ give the difference between the projections $u_{M_i|\mathcal{R}_I}$ and the measured detections $\tilde{u}_{M_i|\mathcal{R}_I}(I_j)$. Therefore, if at least two markers are visible, then the reprojection error of the composite marker on an image I_j can be defined as:

$$E_L(\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_L}) = \sum_{i=1}^{n_m} \sum_{m=1}^4 \sum_{l=1}^2 \left(\Delta u_{M_i|\mathcal{R}_I}(\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_L}) \right)_{l,m,i}^2 \tag{21}$$

By optimization, the locator pose on this image is calculated by minimizing this error by a transformation $\gamma^*(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_L}$ such as:

$$\gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_L} = \arg \min_{\gamma^*(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_L} \in \mathbb{R}^3 \times \mathbb{R}^3} E_L(\gamma^*(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_L}) \tag{22}$$

3 Experiments and results

This section presents the experimental validation of the proposed method. A first part is focused on the calibration and measurement uncertainties of locators made from two different shapes of the host objects. A second one concerns the use of these locators as stylus where their respective tip is tracked. The last part is dedicated to the surface registration on an artificial scapula with points digitized by one of the presented styli. All tests were performed on a laptop using an Intel Core i9-9880H processor with 32 GB of RAM and Ubuntu 18.04 as the operating system. An SVS-VISTEK EXO542 MU3 camera with a resolution of 5320×3032 pixels, 23 fps max and a 16mm focal length PENTAX TV lens were used. The scene is lit by a ring light. The intrinsic camera parameters were identified by the chessboard corner method [41]. An overview of the experimental setup is provided in Fig. 4a. The working distance (distance camera-object) is in the range of 800–1200 mm. The aperture of the

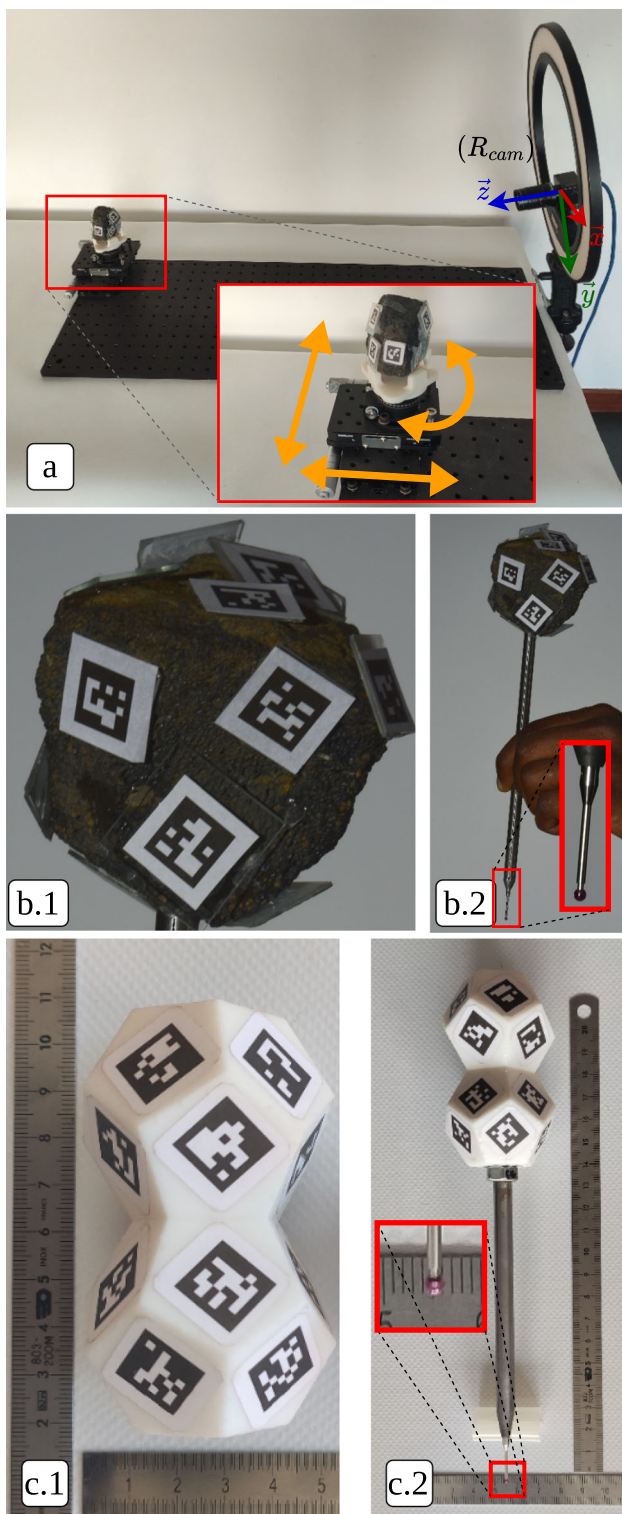


Fig. 4 Experimental setup equipments. **a** Camera, ring light and micrometric stage. The camera reference frame is represented. The micrometric stage allows to perform rotations and translations close to camera x and z axis. **b.1** Locator made from stone, ArUco were glued on glass. **b.2** Stylus made from stoned-based locator. **c.1** The double-dodecahedron used had the dimensions $9\text{ cm} \times 7\text{ cm} \times 7\text{ cm}$. ArUco markers glued on 21 faces of the solid. **c.2** Stylus using the double-decahedron. The spherical contact probe is shown

lens diaphragm is set to obtain a depth of field of about 400 mm.

3.1 Calibration and measurement uncertainties of locators

The first proposed locator is composed by 12mm sided 6×6 ArUco markers placed on the surface of a stone of approximate dimension of $70 \times 70 \times 60\text{ mm}^3$ with a random shape, see Fig. 4b.1. The ArUco markers are stuck on glass plates to guarantee their flatness, themselves glued on the stone. These markers were placed randomly on the surface of the object, simply avoiding placing them in a coplanar orientation. Depending on the geometry of the host object, the number of markers visible to the camera can vary. We propose to analyze another geometry in the form of a double-dodecahedron, see Fig. 4c.1. This geometry offers the advantage of presenting a minimum of 6 ArUco markers whatever the camera's observation point.

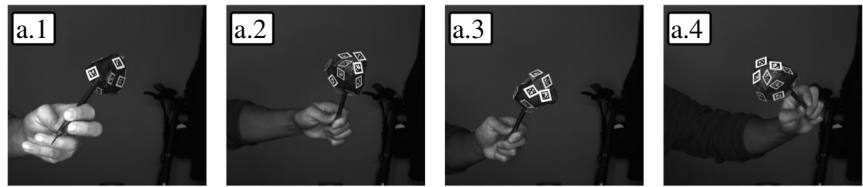
As shown in Fig. 5a, a continuous recording during about 20s of freehand manipulations of the stylus is performed with great variety of poses. This provided a set of 423 frames. Attention was paid to avoid motion blur, overexposure and to respect the camera depth of field.

The proposed algorithm is applied to calibrate the locator made from a stone. The structure of the locator is calculated from the individual poses of each of the ArUco markers on the set of images (2525 detections lead to 2×2525 poses due to the pose ambiguity). Figure 5b represents the graph which connects all the images to the 15 markers glued on the locator. Figure 5c shows the corner positions detected by the AprilTag 2 fiducial marker detection method [42]. The position of these same corners determined by projection with our method alone and after the bundle adjustment are represented jointly on these four images. Some marker's detections, although visible on some images, are excluded by our method because they do not reach the criteria defined by Eqs. 16 and 17. Figure 5.d presents the starting point given by our method and the optimized state. The mean and standard deviation values of these reprojection errors along x and y axes are, respectively: $\bar{X}_x = 3.73e-4$ pixel, $\sigma_x = 1.55e-1$ pixel; $\bar{X}_y = -1.36e-3$ pixel, $\sigma_y = 1.61e-1$ pixel. To sum up, 98% of the points have a reprojection errors contained in a square box of 1 pixel side. The points outside this area are due to markers strongly tilted to the camera's optical axis (approximately greater than 60°). Keeping or removing them only marginally affects the result of the optimization. In practice, the calibration time varies between 3 and 6 min depending on the number of images recorded and the number of present markers.

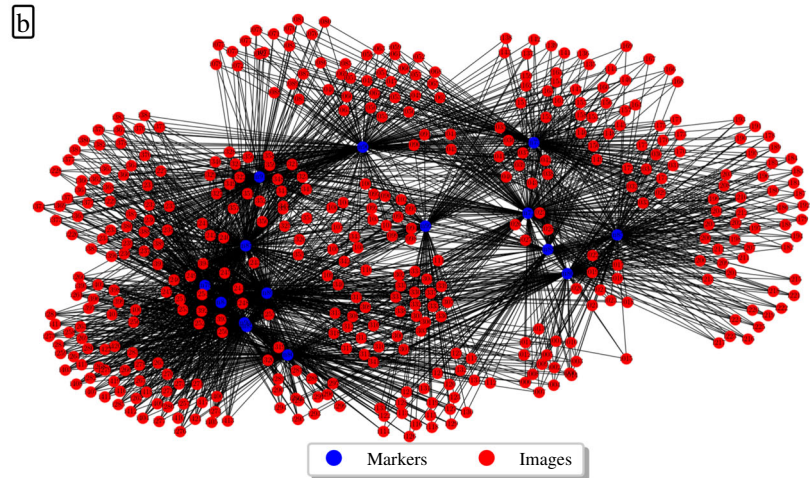
Using Eq. 22, the pose of the locator obtained after calibration can be calculated from an image in lower than 0.1 second, examples are shown in Fig. 6. We propose to qualify

Fig. 5 Workflow for the calibration of a rigid body.
a Sample of four calibration images of the locator object. The stylus was moved manually in front of the camera during image acquisition. All faces of the stone were presented so that all markers are considered during calibration. **b** Initial graph obtained during the calibration of the locator object, red and blue nodes denote images and markers, respectively. The edges represent the detections. **c** Positions of each corner of each marker on the images. In red, the positions obtained using the AprilTag 2 fiduciary detection method [42]. In blue, the positions obtained after projection of the poses estimated by our method. And finally, in green, the positions obtained after the least-squares optimization use the blue positions as initial solutions. **d** Residuals of reprojection errors of the markers corners from initial solution (blue) and the results of least-square optimization (green) on the all image batch. The majority of the points (98%) after optimization presents an absolute value reprojection error lower than 0.5 pixel

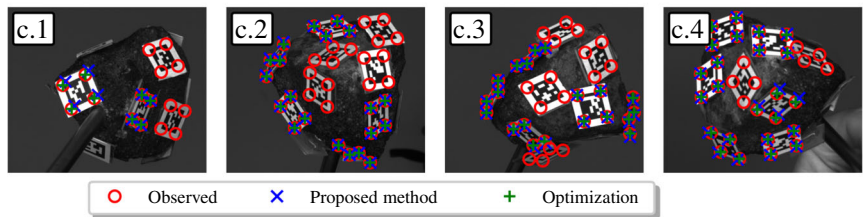
Image batch sample



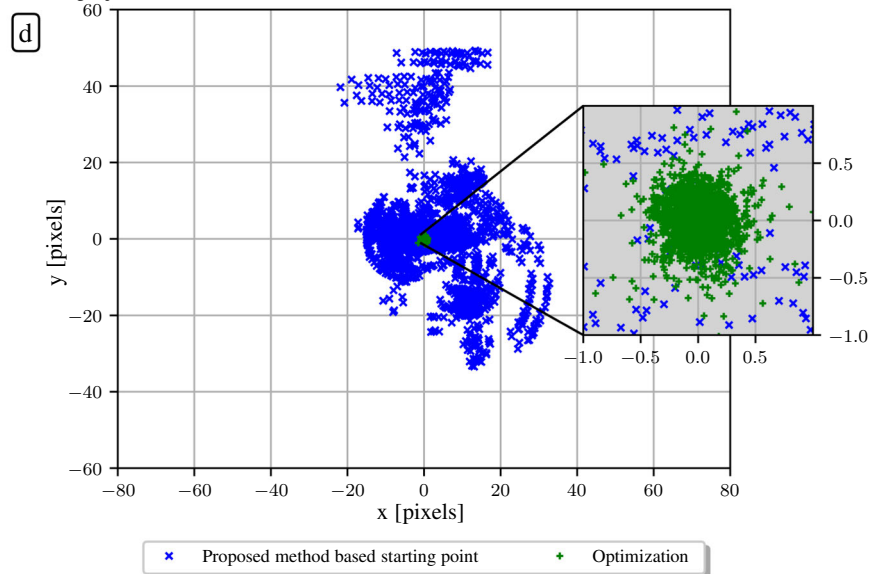
Initial graph



Corner positions



Corner reprojection residues



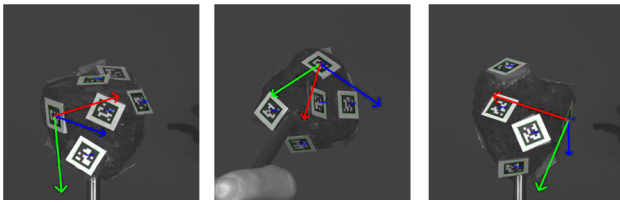


Fig. 6 Different poses of the locator. The pose of the locator can be calculated in real time, regardless of its orientation, if a minimum of two markers are visible. The reference frame associated to the locator corresponds to center of the chosen reference marker M_0 . Note that even if M_0 is not visible on an image, the pose of the locator can be still calculated and drawn using the other markers detections

the measurement uncertainties in translation and in rotation of the resulting locator.

→ Translation uncertainties

From a reference position noted Q_0 , two successive translations of target distance $L_t = 25$ mm are performed on the locator using a micrometric translation stage. The directions of these two translations are orthogonal to each other. The first translation moves the locator from the position Q_0 to position Q_1 along a direction close to the camera optical axis (z axis). The second translation brings the locator from Q_1 to Q_2 with a direction close to the camera x axis. For each position, $N_i = 10$ images are recorded. We calculate the associated Euclidean distance noted L between the different positions. Note that the distances are calculated from all the combinations of images between the two positions, so there are $N_i^2 = 100$ distances calculated between each position. The whole procedure is repeated for 3 orientations of the locator and the associated distances are noted L_l , with $l \in \{0, 1, 2\}$. We impose that the visible markers between the three orientations are different.

The results of these tests are shown in Fig. 7. For all the movements made along the camera pseudo z axis, an average distance of 25.08 mm is recorded. The associated standard deviation dispersion is 0.06 mm. Similarly, the measurements taken on the pseudo x camera axis give an average of 25.01 mm and a standard deviation dispersion of 0.03 mm. According to these measurements, a positioning error of less than one millimeter is possible. As expected, the standard deviation dispersion is higher for the movements along the camera optical than the orthogonal direction.

→ Rotation uncertainties

A qualification of the measurement uncertainties in rotation of the locator is proposed. From an initial orientation, successive imposed rotations of $\theta_t = 45^\circ$ are performed on the locator using a rotation plate. A total of $N_r = 7$ rotations are thus imposed. For each orientation, $N_i = 10$ images are taken. Between each rotation, the pose of the locator is estimated and the value of the rotation angle θ performed is then deduced.

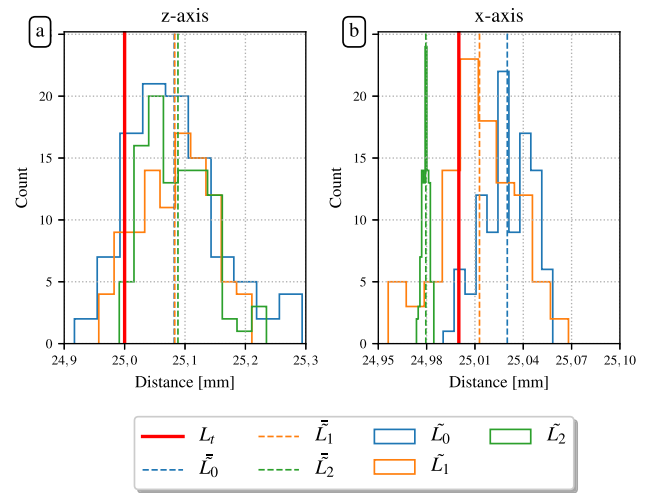


Fig. 7 Translation uncertainties of the locator for $l \in \{0, 1, 2\}$ given orientations. **a** Distribution of the distances L_l calculated between positions Q_0 and Q_1 , close to camera z-axis direction. **b** Distribution of the distances L_l calculated between positions Q_1 and Q_2 , close to camera x axis direction

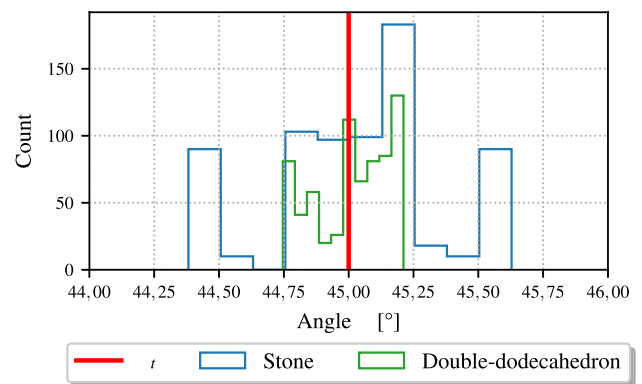


Fig. 8 Rotational uncertainties of the stone locator. From a initial orientation, 7 rotations of $\theta_t = 45^\circ$ are imposed

The results of these tests are given in Fig. 8. For the imposed angles θ , a standard deviation of $\sigma_\theta = 0.33^\circ$ and $\sigma_\theta = 0.14^\circ$, respectively, for the stone and the double-dodecahedron locators.

We applied the same translations and rotations tests to the double-dodecahedron locator geometry. The test results for translational uncertainties along the camera z-axis and rotational uncertainties from the two types of locator geometries are given in Table 2.

There is a significant improvement in the results with the use of the double-dodecahedron, since the number of visible markers is systematically greater.

Fig. 9 Stylus tip identification. **a** Sample calibration images of the locator tip. **b** Illustration of the movements performed during image capture, the tip of the stylus is held in a conical notch. The operator manipulates the stylus while keeping contact with the notch. At the same time, the operator rotates the stylus along the axis of the stainless steel rod to present all the faces of the locator object. The range of motion of the locator is about ± 100 mm around the vertical position. **c** Photography of the tip of the stylus. The origin of the drawn reference frame is positioned on the tip of the stylus

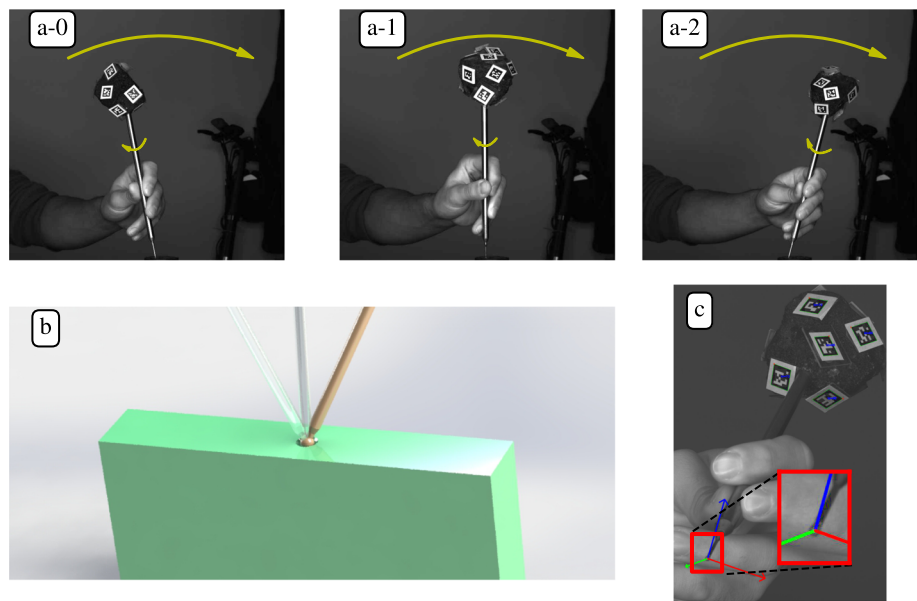


Table 2 Accuracy uncertainties in translation and rotation of the double-dodecahedron shape locator compared to the locator made from a stone

Host object	$ \bar{d}_z - d_t $	σ_{d_z}	d_t^{**}	$ \bar{\theta} - \theta_t $	σ_θ	θ_t^{**}
Stone	0.080	0.060	25	0.041	0.330	45
Dd*	0.002	0.041			0.010	0.140
Unit		[mm]			[°]	

*Double-dodecahedron, **Target values

3.2 Stylus tracking

A steel rod of length 200mm with a spherical contact probe ($\varnothing 2$ mm) was embedded in the stone and the double-dodecahedron locators see Fig. 4b-c.2.

The position of the stylus tip was identified by using movements in which the tip was kept in contact with a fixed conical surface relative to the camera. The performed movements are random combinations of rotations including rotations of 360° around the axis of revolution of the stylus to ensure observation of all markers. Twenty images are recorded as illustrated by Fig. 9.

The center of the tip is determined by looking for a unique point P_{A/\mathcal{R}_L} in the locator frame of reference \mathcal{R}_L that would also be stationary in the camera’s reference frame \mathcal{R}_c during the movements performed. This condition can be expressed as the minimization of the variance of the coordinates of point P_{A/\mathcal{R}_c} in the camera frame:

$$P_{A/\mathcal{R}_L} = \underset{P_{A/\mathcal{R}_L} \in \mathbb{R}^3}{\operatorname{arg\,min}} \left(\operatorname{Var} \left(\left\{ \gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_L} \cdot P_{A/\mathcal{R}_L}^* \right\} \right) \right) \tag{23}$$

The coordinates of stylus tip’s center are thus known in the locator reference frame.

To quantify uncertainties in the position of the tip of this stylus, we reuse the same types of movements as those applied when the tip identification. These motions are performed via a micrometric translation stage for three positions Q_0 , Q_1 and Q_2 of the tip relatively to the camera. The displacement between Q_0 and Q_1 corresponds to a translation of $d_t = 10$ mm in a direction close to the optical axis camera (z axis). The second displacement Q_1 to Q_2 corresponds to a translation of in a direction close to the x axis camera with the same d_t value. For each position, 100 images are recorded.

In order to identify the stability of the process, the standards deviations of the tip coordinates are calculated for the Q_0 position. The results are shown in the following Table 3.

Then, we calculate the distances using the Euclidean norm separating the positions Q_0 and Q_1 noted d_{z_i} and the positions Q_1 and Q_2 noted d_{x_j} . The results are shown in Fig. 10 for the stone locator.

It is thus possible to identify by monocular vision the position of the tip of a stylus observed at a distance of 1 m with a standard deviation of less than 0.2 mm and a mean deviation of less than 0.1 mm.

The comparison of the uncertainties related to the two styli is given in table Tab. 4.

3.3 Surface registration on an artificial scapula

In the context of surgical navigation, we need to identify the position of a bone structure relative to the camera by simply observing a locator rigidly attached to the bone structure. To illustrate our approach, we propose to find the position of

Table 3 Comparison of the stability of tip point identification for the stone and the double-dodecahedron locators

Host object	σ_x	σ_y	σ_z
Stone	0.120	0.040	0.240
Dd*	0.096	0.039	0.302
Unit	[mm]		

The results are obtained by calculating the standard deviations of the coordinates of the tip point in the camera frame. *Double-dodecahedron

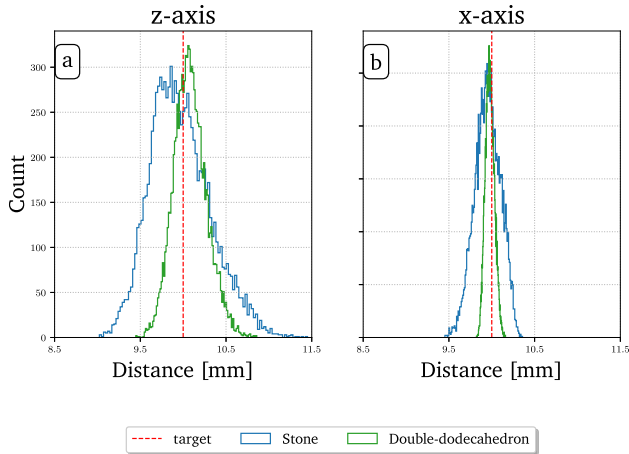


Fig. 10 Styli-tip distance of the locators dispersion where target value is 10 mm. **a)** Distribution of distances d_{z_i} measured between positions P_{0/\mathcal{R}_c} and P_{1/\mathcal{R}_c} along an axis close to the z-axis camera. **b)** Distribution of distances d_{x_i} measured between positions P_{1/\mathcal{R}_c} and P_{2/\mathcal{R}_c} along an axis close to the x-axis camera

Table 4 Comparison of the uncertainties on the position of the tip of the styli between the localizers made from a stone and a double-decahedron

Host object	$ \bar{d}_{z_i} - d_t $	$\sigma_{d_{z_i}}$	$ \bar{d}_{x_i} - d_t $	$\sigma_{d_{x_i}}$	d_t^{**}
Stone	0.063	0.344	0.058	0.169	10
Dd*	0.081	0.193	0.022	0.053	
Unit	[mm]				

*Double-dodecahedron, **Target distance value

an artificial scapula in relation to the camera using a surface registration. The shape of the used locators were double-dodecahedron geometry.

In the example presented here, a digital scapula model (Fig. 11a) was obtained from a CT scan of a patient. This STL model was used to reproduce an artificial copy of the scapula by photopolymerization using the *Form 3+* printer from the manufacturer *Formlabs*. In this test, we propose to use two locators, one attached to the stylus and the second attached to the artificial scapula whom references frames are, respectively, \mathcal{R}_{sty} and \mathcal{R}_{scp} . The coordinates of the points of the numerical model of the scapula are known in a reference frame noted \mathcal{R}_{stl} . Using the stylus, scans of certain areas of

the artificial scapula are made during which $N_i = 90$ images are taken (Fig. 11b).

For each image $j \in \{1, \dots, N_i\}$, the center of the spherical tip of the stylus noted $P_{C_j/\mathcal{R}_{scp}}$ is known in the reference frame \mathcal{R}_{scp} by the following equation:

$$P_{C_j/\mathcal{R}_{scp}} = \gamma(I_j)_{\mathcal{R}_{scp} \leftarrow \mathcal{R}_c} \cdot \gamma(I_j)_{\mathcal{R}_c \leftarrow \mathcal{R}_{sty}} \cdot P_{C/\mathcal{R}_{sty}} \tag{24}$$

Using the proximity.signed_distance() f function from *Trimesh* module (Python3), we can calculate the signed distance, noted Is , between any point expressed in \mathcal{R}_{stl} reference and the surface S forming the numerical model. The sign of Is is negative when the considered point is outside the surface S .

During scanning, the actual probe point is shifted by the radius of the contact sphere of radius $r = 1$ mm. Taking into account the offset, we can write the following optimization equation to identify the rigid transformation $\gamma_{\mathcal{R}_{stl} \leftarrow \mathcal{R}_{scp}}$ between the digital model and the scapula locator as:

$$\gamma_{\mathcal{R}_{scp} \leftarrow \mathcal{R}_{stl}} = \arg \min_{\gamma_{\mathcal{R}_{scp} \leftarrow \mathcal{R}_{stl}}} \sum_{j=1}^{N_i} \left[f \left(\gamma_{\mathcal{R}_{stl} \leftarrow \mathcal{R}_{scp}}^* \cdot P_{C_j/\mathcal{R}_{scp}}, S \right) + r \right]^2 \tag{25}$$

Then, the coordinates of the points $P_{C_j/\mathcal{R}_{stl}}$ can then be expressed in the reference frame \mathcal{R}_{stl} according to the following equation:

$$P_{C_j/\mathcal{R}_{stl}} = \gamma_{\mathcal{R}_{stl} \leftarrow \mathcal{R}_{scp}} \cdot P_{C_j/\mathcal{R}_{scp}} \tag{26}$$

We have repositioned these points in the STL reference frame to represent them on the surface of the numerical model (Fig. 11c). We evaluated the solution found by calculating the obtained signed distance Is between the points $P_{C_j/\mathcal{R}_{stl}}$ and the surface S forming the numerical model.

The results obtained are presented in Fig. 12. The mean value of the signed distance $Is_j + r = 0.011$ mm is observed with a standard deviation $\sigma_{Is_j} = 0.26$ mm.

4 Discussion

Based on the aforementioned results, we evaluate the quality of the following aspects: calibration of locator, the enhanced value associated with the geometry of the double-dodecahedron, positioning of the tip of the stylus, and surface registration for the presented medical application.

A calibration step allows to find the relative positions of the planar markers stuck on these objects by simply observing images of them manipulated freehand in front of the camera. This calibration method is based on graph theory. After

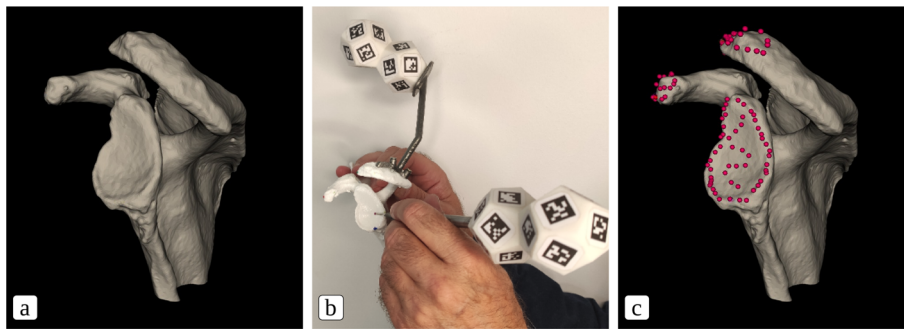


Fig. 11 **a** Numerical model of the scapula from CT scan of cadaveric scapula. **b** Scanning of the artificial scapula. One of the double-dodecahedron shape locators is rigidly attached to the artificial scapula. The scapula and pen are manipulated freehand. The digitalization of

the points is made using the stylus with a second double-dodecahedron shape locator. **c** Registration of the scanned points on the artificial scapula with the digital model of the scapula

optimization, we found the structure of the locator. The reprojection errors observed on the images were mostly contained in a square of 1 pixel.

To test the found solution, controlled movements were imposed on the locators alone. The analysis of the uncertainties for these movements showed that the geometry of double-dodecahedron with at least 6 planar markers performs best. For this locator, we obtain uncertainties of the order of $41\mu\text{m}$ in translation along the camera optical axis (the most critical axis in monocular vision) and uncertainties in rotation of 0.14° .

For the specific needs of shoulder surgery, a probing device is necessary. A stylus of approximately 20 cm in length consisting of a spherical tip and a double-dodecahedron locator is proposed. Tests from paragraph 3.2 have shown that it is possible to track the center of this sphere with an uncertainty of about 0.2 mm.

This stylus was used to perform a surface registration between an artificial scapula and its digital model based on the scanning of a restricted area of this scapula (glenoid cavity). The deviations between the set of scanned points and the mesh of the digital model are all less than 1 mm. This accuracy is fully compatible with the expectations of orthopedic shoulder surgery. The results obtained from cadaveric shoulders will be presented in a forthcoming publication. The calibration time of the locators and the identification of the stylus tip do not impact the operating time since these steps can be performed preoperatively. However, the palpation of the glenoid cavity and the surface registration are included in the operating time (about 2 min in total).

Improvement perspectives are possible for this monocular localization. Currently, the calibration of the locator with a hundred images in input takes about 1 min. The cycle basis obtained from the graph tool proposes short cycles but also longer cycles, the latter penalizing the processing time. A reflection to limit this type of long cycle has to be conducted.

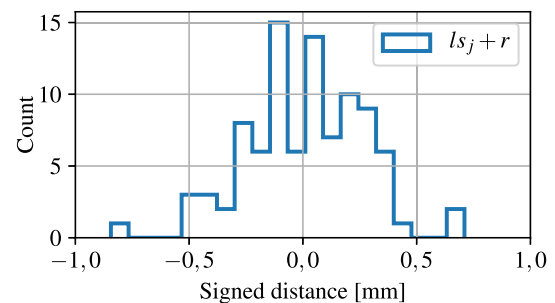


Fig. 12 Distribution of the signed distance $ls_j + r$ between the measured points $P_{C_j/\mathcal{R}_{st}}$ and the surface S of the digital scapula model taking into account the offset of the contact sphere radius r

More generally, structural optimizations of the algorithm should allow an improvement of the processing speed using parallelization approaches.

5 Conclusion

In this work, a method to track in real time an object positioned at about 1 m from the camera by monocular vision was proposed. From a random shape object (stone) or from a more controlled geometry (double-dodecahedron), 3D locators based on a cluster of planar markers were analyzed.

The resulting locators have translation uncertainties on the order of one-tenth of a millimeter and orientation uncertainties on the order of one-hundredth of a degree. We have also shown that these locators can be used for surface registration with sub-millimeter uncertainties. These uncertainties are reached whatever the visible plane markers of this locator. These results, along with the locators' geometric versatility, show great potential for their implementation in orthopedic surgical navigation applications. In this application, locators can be customized for use with medical tools such as probing styluses, drills, surgical saws, implants or prostheses.

The freedom offered by this approach allows to imagine applications for other types of surgeries but also for various fields of applications requiring real-time sub-millimetric localizations between different objects.

Acknowledgements This work was financed by the European program INTERREG France-Switzerland.

References

- Fischler, M.A., Bolles, R.C.: Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM* **24**(6), 381–395 (1981)
- Li, S., Xu, C., Xie, M.: A robust $O(n)$ solution to the perspective- n -point problem. *IEEE Trans. Pattern Anal. Mach. Intell.* **34**(7), 1444–1450 (2012)
- Schweighofer, G., Pinz, A.: Robust pose estimation from a planar target. *IEEE Trans. Pattern Anal. Machine Intell.* **28**, 2024–2030 (2006)
- Wang, J., Sadler, C., Montoya, C.F., Liu, J.C.: Optimizing ground vehicle tracking using unmanned aerial vehicle and embedded apriltag design. In: 2016 International conference on computational science and computational intelligence (CSCI), pp. 739–744, (2016)
- Sani, M.F., Karimian, G.: Automatic navigation and landing of an indoor AR. drone quadrotor using ArUco marker and inertial sensors, In: 2017 International conference on computer and drone applications (ICoNDA), (Kuching), pp. 102–107, IEEE, (2017)
- Ishiyama, H., Kurabayashi, S.: Monochrome glove: a robust real-time hand gesture recognition method by using a fabric glove with design of structured markers, In: 2016 IEEE virtual reality (VR), pp. 187–188, ISSN: 2375-5334 (2016)
- Zhang, X., Fronz, S., Navab, N.: Visual marker detection and decoding in AR systems: a comparative study, In: Proceedings international symposium on mixed and augmented reality, pp. 97–106, (2002)
- Oberkampff, D., DeMenthon, D.F., Davis, L.S.: Iterative pose estimation using coplanar feature points. *Comput. Vis. Image Underst.* **63**, 495–511 (1996)
- Newcombe, R.A., Lovegrove, S.J., Davison, A.J.: Dtam: dense tracking and mapping in real-time, In: 2011 international conference on computer vision, pp. 2320–2327, IEEE, (2011)
- Kendall, A., Cipolla, R.: Modelling uncertainty in deep learning for camera relocation, In: 2016 IEEE international conference on Robotics and Automation (ICRA), pp. 4762–4769, IEEE, (2016)
- Hartley, R., Zisserman, A.: Multiple view geometry in computer vision. Cambridge University Press, Cambridge (2003)
- Snavely, N., Seitz, S.M., Szeliski, R.: Photo tourism: exploring photo collections in 3d, In: ACM siggraph 2006 papers, pp. 835–846, (2006)
- Newcombe, R.A., Izadi, S., Hilliges, O., Molyneaux, D., Kim, D., Davison, A.J., Kohi, P., Shotton, J., Hodges, S., Fitzgibbon, A.: Kinectfusion: real-time dense surface mapping and tracking, In: 2011 10th IEEE international symposium on mixed and augmented reality, pp. 127–136, Ieee, (2011)
- Mur-Artal, R., Tardós, J.D.: Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras. *IEEE Trans. Rob.* **33**(5), 1255–1262 (2017)
- Moré, J.J.: The levenberg-marquardt algorithm: implementation and theory, In Numerical analysis: proceedings of the biennial conference held at Dundee, June 28–July 1, 1977, pp. 105–116, Springer, (2006)
- Levenberg, K.: A method for the solution of certain non-linear problems in least squares. *Q. Appl. Math.* **2**(2), 164–168 (1944)
- Marquardt, D.W.: An algorithm for least-squares estimation of non-linear parameters. *J. Soc. Ind. Appl. Math.* **11**(2), 431–441 (1963)
- Triggs, B., McLauchlan, P.F., Hartley, R.I., Fitzgibbon, A.W.: Bundle adjustment—a modern synthesis, In: Vision algorithms: theory and practice: international workshop on vision algorithms Corfu, Greece, September 21–22, 1999 Proceedings, pp. 298–372, Springer, (2000)
- Lepetit, V., Moreno-Noguer, F., Fua, P.: Ep n p: an accurate $O(n)$ solution to the p n p problem. *Int. J. Comput. Vision* **81**, 155–166 (2009)
- Lourakis, M., Terzakis, G.: A globally optimal method for the pnp problem with mrp rotation parameterization, In: 2020 25th international conference on pattern recognition (ICPR), p. 3058–3063, (2021)
- Nister, D.: An efficient solution to the five-point relative pose problem. *IEEE Trans. Pattern Anal. Mach. Intell.* **26**, 756–770 (2004)
- Barath, D., Polic, M., Förstner, W., Sattler, T., Kukulova, Z.: Making a ngeecoomrertersypocnodmenpcuetsawtioonrk in camera,” p. 17
- Ségvic, S., Schweighofer, G., Pinz, A.: Performance evaluation of the five-point relative pose with emphasis on planar scenes1, p. 8
- Muñoz-Salinas, R., Marín-Jimenez, M.J., Medina-Carnicer, R.: SPM-SLAM: simultaneous localization and mapping with squared planar markers. *Pattern Recogn.* **86**, 156–171 (2019)
- Garrido-Jurado, S., Muñoz-Salinas, R., Madrid-Cuevas, F., Marín-Jiménez, M.: Automatic generation and detection of highly reliable fiducial markers under occlusion. *Pattern Recogn.* **47**, 2280–2292 (2014)
- Wu, P.-C., Wang, R., Kin, K., Twigg, C., Han, S., Yang, M.-H., Chien, S.-Y.: DodecaPen: accurate 6DoF tracking of a passive stylus, In: Proceedings of the 30th annual ACM symposium on user interface software and technology, UIST '17, (Québec City, QC, Canada), pp. 365–374, Association for Computing Machinery, (2017)
- Wu, C.: Towards linear-time incremental structure from motion, In: 2013 International conference on 3D vision-3DV 2013, pp. 127–134, IEEE, (2013)
- Triggs, B., McLauchlan, P.F., Hartley, R.I., Fitzgibbon, A.W.: Bundle Adjustment—a modern synthesis. In: Goos, G., Hartmanis, J., van Leeuwen, J., Triggs, B., Zisserman, A., Szeliski, R. (eds.) Vision algorithms: theory and practice, pp. 298–372. Springer, Berlin Heidelberg (2000)
- Schonberger, J.L., Frahm, J.-M.: Structure-From-Motion Revisited, pp. 4104–4113, (2016)
- Kotake, D., Uchiyama, S., Yamamoto, H.: A marker calibration method utilizing a priori knowledge on marker arrangement, In: Third IEEE and ACM international symposium on mixed and augmented reality, pp. 89–98, (2004)
- Kato, H., Billinghurst, M., Poupyrev, I., Imamoto, K., Tachibana, K.: Virtual object manipulation on a table-top ar environment, In: ISAR 2000: proceedings of the IEEE and ACM international symposium on augmented reality, pp. 111–119, (2000)
- Munoz-Salinas, R.: Aruco: a minimal library for augmented reality applications based on opencv,” Universidad de Córdoba, **386**, (2012)
- Olson, E.: Apriltag: A robust and flexible visual fiducial system, In: 2011 IEEE international conference on robotics and automation, pp. 3400–3407, IEEE, (2011)
- Fiala, M.: ARTag, a fiducial marker system using digital techniques, In: 2005 IEEE computer society conference on computer vision and pattern recognition (CVPR '05), vol. 2, (San Diego, CA, USA), pp. 590–596, IEEE, (2005)
- Romero-Ramirez, F.J., Muñoz-Salinas, R., Medina-Carnicer, R.: Speeded up detection of squared fiducial markers. *Image Vis. Comput.* **76**, 38–47 (2018)

36. Garrido-Jurado, S., Muñoz-Salinas, R., Madrid-Cuevas, F.J., Medina-Carnicer, R.: Generation of fiducial marker dictionaries using mixed integer linear programming. *Pattern Recogn.* **51**, 481–491 (2016)
37. Wang, J., Olson, E.: AprilTag 2: efficient and robust fiducial detection, In: 2016 IEEE/RSJ international conference on intelligent robots and systems (IROS), pp. 4193–4198, ISSN: 2153-0866 (2016)
38. Collins, T., Bartoli, A.: Infinitesimal plane-based pose estimation. *Int. J. Comput. Vis.* **109**, 252–286 (2014)
39. Paton, K.: An algorithm for finding a fundamental set of cycles of a graph. *Commun. ACM* **12**, 514–518 (1969)
40. Dijkstra, E.W., et al.: A note on two problems in connexion with graphs. *Numer. Math.* **1**(1), 269–271 (1959)
41. De la Escalera, A., Armingol, J.M.: Automatic chessboard detection for intrinsic and extrinsic camera parameter calibration. *Sensors* **10**(3), 2027–2044 (2010)
42. Wang, J., Olson, E.: Apriltag 2: efficient and robust fiducial detection, In: 2016 IEEE/RSJ international conference on intelligent robots and systems (IROS), pp. 4193–4198, IEEE, (2016)

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.