



Tracking objects with generic calibrated sensors: An algorithm based on color and 3D shape features^{☆,☆☆}

M. Taiana^{*}, J. Santos, J. Gaspar, J. Nascimento, A. Bernardino, P. Lima

Institute for Systems and Robotics, Instituto Superior Técnico, 1049-001 Lisboa, Portugal

ARTICLE INFO

Article history:

Available online 11 March 2010

Keywords:

Omnidirectional vision
Tracking
Particle filtering

ABSTRACT

We present a color and shape based 3D tracking system suited to a large class of vision sensors. The method is applicable, in principle, to any known calibrated projection model. The tracking architecture is based on particle filtering methods where each particle represents the 3D state of the object, rather than its state in the image, therefore overcoming the nonlinearity caused by the projection model. This allows the use of realistic 3D motion models and easy incorporation of self-motion measurements. All nonlinearities are concentrated in the observation model so that each particle projects a few tens of special points onto the image, on (and around) the 3D object's surface. The likelihood of each state is then evaluated by comparing the color distributions inside and outside the object's occluding contour. Since only pixel access operations are required, the method does not require the use of image processing routines like edge/feature extraction, color segmentation or 3D reconstruction, which can be sensitive to motion blur and optical distortions typical in applications of omnidirectional sensors to robotics. We show tracking applications considering different objects (balls, boxes), several projection models (catadioptric, dioptric, perspective) and several challenging scenarios (clutter, occlusion, illumination changes, motion and optical blur). We compare our methodology against a state-of-the-art alternative, both in realistic tracking sequences and with ground truth generated data.

© 2010 Elsevier B.V. All rights reserved.

1. Introduction

Omnidirectional and wide-angle vision sensors have been widely used in the last decade for robotics and surveillance systems. These sensors gather information from a large portion of the surrounding space, thus reducing the number of cameras required to cover a certain spatial extent. Their classical applications include mobile robot self-localization and navigation [1,2], video surveillance [3] and humanoid foveal vision [4]. One drawback is that images suffer strong distortions and perspective effects, demanding non-standard algorithms for target detection and tracking.

In scenarios where the shape of objects can be modeled accurately *a priori*, 3D model based techniques are among the most successful in tracking and pose estimation applications [5]. However,

classical 3D model based tracking methods are strongly dependent on the projection models and, thus, are not easily applicable to omnidirectional images. Often nonlinear optimization methods are employed: a cost function expressing the mismatch between the predicted and observed object points is locally minimized with respect to the object's pose parameters [5]. This process involves the linearization of the relation between state and measurements, which can be very complex with omnidirectional sensor geometries. These approaches have limited convergence basins requiring either small target motions or very precise prediction models.

In this paper we try to overcome these problems by addressing the pose estimation and tracking problem in a Monte Carlo sampling framework [6], through the use of a Particle Filter (PF).

Particle Filters (PF) have become a popular tool in image processing and computer vision communities as a way to infer time varying properties of a scene from a set of input images. PF compute a sampled representation of the posterior probability distribution over the scene properties of interest. It is capable of coping with nonlinear dynamics and nonlinear observation equations, which helps in easily maintaining uncertainty with multiple distributions.

The principle of the PF is simple. The goal is to compute the posterior probability distribution $p(x_t | y_{1:t})$ over an unknown state x_t , conditioned on image observations up to that time instant, i.e., $y_{1:t}$. Particle filtering works by approximating the posterior

[☆] Parts of this manuscript were previously presented at the British Machine Vision Conference (BMVC'08) and at the Workshop on Omnidirectional Robot Vision, held in conjunction with the SIMPAR 2008 conference.

^{☆☆} This work was supported by the European Commission, Projects IST-004370 RobotCub and ICT-231640 HANDLE, and by the Portuguese Government - Fundação para a Ciência e Tecnologia (ISR/IST pluriannual funding) through the PIDDAC program funds, and through project BIO-LOOK, PTDC/EEA-ACR/71032/2006.

^{*} Corresponding author. Tel.: +351 960323163.

E-mail addresses: mtaiana@isr.ist.utl.pt, matteo.taiana@gmail.com (M. Taiana), jsantos@isr.ist.utl.pt (J. Santos), jag@isr.ist.utl.pt (J. Gaspar), jan@isr.ist.utl.pt (J. Nascimento), alex@isr.ist.utl.pt (A. Bernardino), pal@isr.ist.utl.pt (P. Lima).

density using a discrete set of samples, *i.e.*, the states. Each state corresponds to some hypothesized set of model parameters. Each sample is typically weighted by its likelihood $p(y_t|x_t)$, the probability that the current state observations were generated by the hypothesized state. Each sample can be considered as a state hypothesis, whose weight depends on the corresponding image data. The method operates by “testing” state hypothesis, thus avoiding the linearization between state and measurements required in gradient optimization techniques. In the context of our problem, this allows the utilization of arbitrarily complex projection models.

More precisely, 3D object localization hypothesis, generated by the particle filter, allow obtaining the 2D appearance of the objects, provided one has the projection models of the imaging sensors. Examples of imaging sensors used in our experiments range from conventional perspective to lens-mirror (catadioptric) and fisheye-lens based omnidirectional cameras. Note that while the perspective projection is already a nonlinear model, in the sense that constant 3D motion increments of an object do not imply constant 2D image increments, it is still a simple model in the sense that the projection of any 3D point is ray-traced as a straight line passing through one single point. This is not true in general for fisheye or catadioptric cameras. Fisheye lenses bend the incoming principal optical rays progressively more while moving towards the periphery of the field of view. Catadioptric cameras use simple lenses, but the mirror implies a reflection according to the its local slope. In all cases one can have strong, anisotropic, geometrical distortions associated to the projection which bend and blur the object’s visual features in a space variant manner, creating difficulties in the search for object’s features. In our approach, adopting a 3D tracking context, we can use the projection model to predict and test the location of features in the images instead of locally searching for them.

1.1. Related work

Particle filtering methods have been extensively used during the last decade in 2D tracking applications. One of the first algorithms applying particle filters in the 2D context was the Condensation algorithm [7]. In that work targets were modeled with a contour based appearance template. The approach proved not very robust to occlusion and hard to initialize. To address such limitations, more recent works added other types of features such as color histograms [8] and subspace methods [9]. In [10] a hierarchical particle filter is used to track objects (persons) described by color rectangle features and edge orientation histogram features. Since the use of multiple cues increases the computational demands, two optimized algorithmic techniques are employed: integral images are used to speed-up feature computation and an efficient cascade scheme is used to speed-up particle likelihood evaluation. To deal with occlusions [11] fuses color histogram features acquired from multiple cameras. Cameras are pre-calibrated with affine transformations and the state of the filter is composed by the 2D coordinates and bounding box dimensions of the target in one of the cameras. Each particle in the filter generates a bounding box on each camera and the observation model is composed by the concatenation of the color histograms in all bounding boxes. It is shown that multiple cameras’ tracking and data fusion are able to tackle situations when the target is occluded in one of the cameras.

Despite the success of particle filters in 2D tracking applications, not many works have proposed their use in a 3D model based context. In [12] it is proposed a system for estimating the time varying 3D shape of a moving person having 28DOF. To deal with the high dimensionality, hybrid Monte Carlo (HMC) filter is used. However in that work observations were obtained directly from the 2D projection of suitably placed markers in the human body and, therefore is only applicable in restricted cases. In

general settings, the true appearance of an object (e.g. color, shape, contours) must be taken into account. For instance, the work in [13] uses a particle filter [7] to implement a full 3D edge tracking system using fast GPU computation for real-time rendering of each state hypothesis image appearance (visible edges) and for applying edge-detectors in incoming video stream. In [14] the computation of edges in the full image is avoided by grouping line segments from a known model into 3D junctions and forming fast inlier/outlier counts on projected junction branches. A local search for edges around the expected values must be performed at each time step. In [15] a particle filter is used to estimate the 3D positions of humans. The environment is explicitly modeled to handle occlusions caused by fixed objects. Multiple fixed cameras and background subtraction techniques are used for the computation of the likelihood.

1.2. Our approach

Most previous works on 3D model based tracking, both the ones based on nonlinear optimization and the ones relying in sampling methods, require the extraction of edge points from the images, either by processing the full image with edge detectors or performing local search for edge points. We stress that the primary disadvantage of edge based tracking algorithms is the lack of robustness to motion and optical blur. These effects are frequent in robotics applications due to the mobility of the devices and the frequent out-of-focus situations.

We formulate the problem differently. Rather than determining explicitly the location of edge points in the image, we compute differences between color histograms on the inside and outside regions of the projected target surfaces. We do not require explicit image rendering, in the sense of creating an image of the expected appearance of the target as in [13]. Instead we just need to compute the 3D to 2D projection of some selected points inside and outside the target’s visible surfaces. This can be easily done with any single projection center sensor, and allows a fast evaluation of each particle’s likelihood. Because explicit edge or contour extraction is avoided, the method is more robust to blur arising either from fast object motions or optical defocus. Altogether, our approach facilitates the application to arbitrary nonlinear image projections and settings with fast target/robot motion.

In our approach we use color features to compute state likelihoods. The reasons for using color features are the following: color features cope well with motion and optical blur, which are frequent in the scenarios we consider; color features do not require local image processing for extracting edges or corners, but simply pixel evaluations, which makes the system suitable for real-time applications; finally, many robotics research problems assume objects with distinctive colors to facilitate the figure-ground segmentation problem, for instance in cognitive robotics [16] or robotic competitions [17]. Notwithstanding, the approach is general and, in other scenarios, additional features could be used.

In [17] we have presented the first application of our method for tracking balls (spheres) and robots (cylinders) in the RoboCup Middle Size League scenario with catadioptric sensors. In [18] we have compared the application of two well known approximations of the projection function for the class of perspective catadioptric mirrors in our tracking framework: the unified projection model and the perspective model. In [19] we extended the approach to consider not only the 3D position but also the orientation of the targets and presented applications with dioptric and perspective cameras with radial distortion to track both spherical and convex polyhedral shapes. In [20] we have extended the observation model to track objects without a initial color model (only the shape model is used) and have extended the motion model to consider the observer’s self-motion. In the present work we synthesize

the main results derived in previous work and provide a better characterization of the method's performance in challenging scenarios, including the comparison with an alternative state-of-art technique and quantitative evaluations with ground truth data.

The paper is organized as follows: Section 2 presents common imaging systems used in robotics and corresponding projection models. Section 3 describes the particle filtering approach and the state representation in our problem. In Section 4 we detail the 3D shape and color based observation model used in the tracking filter. In Section 5 we show several experiments that illustrate the performance of the system in realistic scenarios, including a comparison to an alternative approach with ground truth data. In Section 6 we draw conclusions and present directions for future work.

2. Imaging systems

We start with a brief introduction to the most common imaging geometries employed in robotics that were used in our experiments. We focus on imaging systems with axial symmetry, both dioptric and catadioptric. Cameras with axial symmetry can be described by a projection function \mathcal{P} :

$$\rho = \mathcal{P}([r \ \varphi \ z]^T; \vartheta) \quad (1)$$

where the z axis coincides with the optical axis, $[r \ \varphi \ z]^T$ represents a 3D point in cylindrical coordinates, ρ is the radial coordinate of the imaged point (the angular coordinate coincides with angular coordinate of the 3D point, φ) and ϑ is a vector of parameters characterizing the geometry of the system (see Fig. 1a and c).

2.1. Dioptric systems

Conventional lens-only systems, i.e., dioptric systems, have been traditionally described by the pin-hole or perspective projection model (PPM):

$$\rho = k \cdot \frac{r}{z} \quad (2)$$

where k includes the camera's intrinsic calibration parameters. This is valid only for narrow fields of view, as otherwise the radial distortion becomes too significant. Modeling super-fisheye fields of view (views larger than 180° angles) involves using other projection models [21,22]. For example [21] uses a super-fisheye Nikon F8 and proposes the model $\rho = a \cdot \tan(\theta/b) + c \cdot \sin(\theta/d)$, where θ denotes the angle of an incoming optical ray with the optical axis (see Fig. 1c) and (a, b, c, d) are parameters characterizing the optical system, obtained by a calibration procedure. Some recent advances allowed companies to build super-fisheye lenses having simple projection models. For example the Sunex's DSL215 lens, Fig. 1d, is designed to have a Constant Angular Resolution Projection Model (CARPM, termed *equidistant* in [21]), meaning that the radial coordinate of the image points ρ is in a linear relationship with the direction of the incoming optical ray, θ :

$$\rho = f \cdot \theta = f \cdot \arctan(r/z). \quad (3)$$

2.2. Catadioptric systems

The conjunction of lens based systems with mirrors, i.e. catadioptric systems, allows the acquisition of wide-angle and omnidirectional images with relatively simple and cheap setups. As compared to dioptric systems, the mirror introduces a reflection on the incoming optical rays implying that the directions θ observed by the camera effectively correspond to (distinct) directions, ϕ , in the optical rays projecting 3D points towards the mirror (see Fig. 1a).

Deriving \mathcal{P} for the complete 3D field of view (FOV) of a catadioptric system involves using the actual mirror shape, F , which is a function of the radial coordinate t . Based on first order optics, particularly on the reflection law at the specular surface of revolution, (t, F) , we have:

$$\arctan(\rho) + 2 \cdot \arctan(F') = \phi \quad (4)$$

where $\arctan(\rho) = \arctan(t/F) = \theta$ is the angle of the principal optical ray with the optical axis of the system, $\phi = -\arctan((r-t)/(z-F))$ is the angle of the principal optical ray before the reflection on the mirror, and F' represents the slope of the mirror shape. If F is a known shape then Eq. (4) describes a generic Catadioptric Projection Model (CPM), as it forms an equation on ρ for a given 3D point (r, z) .

In general there is no closed-form solution for such model. Some simple shapes F , such as a hyperboloid having one focus coinciding with the center of the pin-hole camera capturing the mirror image, allow deriving a closed-form solution for finding ρ from (r, z) . An interesting solution is proposed in [23,24], where mirrors are designed to have a constant resolution, wide-angle view of the ground plane. Replacing $\rho = t/F$ and $\phi = -\arctan((r-t)/(z-F))$, becomes a differential equation, expressing the constant horizontal resolution property, $\rho = a \cdot r + b$, for one plane $z = z_0$. F is usually found using numerical integration methods [23,24].

A closed form model that provides good approximations to the general CPM is given by the Unified Projection Model (UPM) [25]. It represents all omnidirectional cameras with a single center of projection and consists of a two-step mapping via a unit-radius sphere [25]: (i) a 3D world point, $P = [r \ \varphi \ z]^T$, is projected orthogonally to the sphere surface onto a point P_s ; (ii) projected to a point on the image plane, $P_i = [\rho \ \varphi]^T$, from a point O on the vertical axis of the sphere, through the point P_s . The mapping is defined by:

$$\rho = \frac{l+m}{l\sqrt{r^2+z^2}-z} \cdot r \quad (5)$$

where the (l, m) parameters describe the type of camera. The UPM is a widely used representation for CPM when F describes (i) a hyperboloid or ellipsoid with focus at $(0, 0)$; (ii) a paraboloid combined with a telecentric lens ($\theta = 0$) or (iii) $F = \text{const}$ [26, 27]. Note that the PPM is a particular case of the UPM, obtained by setting $l = 0$ and $m = -k$.

Concluding, in the context of our tracking methodology we can use arbitrarily complex projection models. Examples of models used in our experiments comprise the PPM, CARPM, CPM and UPM, described by Eqs. (2)–(5), or by the general expression Eq. (1).

3. 3D tracking with particle filters

We are interested in estimating, at each time step, the 3D pose of the target. Thus, the state vector of the target, denoted as \mathbf{X}_t , contains its 3D pose and derivatives up to a desired order. It represents the object evolution along time, which is assumed to be an unobserved Markov process with some initial distribution $p(\mathbf{x}_0)$ and a transition distribution $p(\mathbf{x}_t|\mathbf{x}_{t-1})$. The observations $\{\mathbf{y}_t; t \in \mathbb{N}\}$, $\mathbf{y}_t \in \mathbb{R}^{n_y}$, are conditionally independent given the process $\{\mathbf{x}_t; t \in \mathbb{N}\}$ with distribution $p(\mathbf{y}_t|\mathbf{x}_t)$, where n_y is the dimension of the observation vector.

In a statistical setting, the problem is posed as the estimation of the *a posteriori* distribution of the state given all observations $p(\mathbf{x}_t|\mathbf{y}_{1:t})$. Under the Markov assumption:

$$p(\mathbf{x}_t|\mathbf{y}_{1:t}) \propto p(\mathbf{y}_t|\mathbf{x}_t) \int p(\mathbf{x}_t|\mathbf{x}_{t-1}) p(\mathbf{x}_{t-1}|\mathbf{y}_{1:t-1}) d\mathbf{x}_{t-1}. \quad (6)$$

The above equation shows that the *a posteriori* distribution can be computed recursively, using the estimate at the preceding time

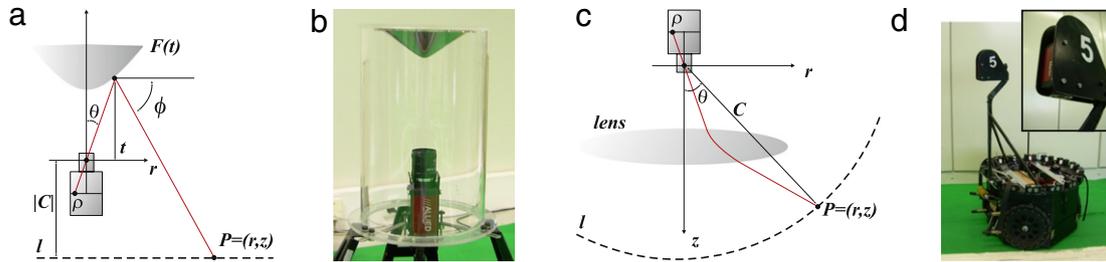


Fig. 1. Catadioptric and dioptric cameras. (a, b) CPM geometry and setup. FOV of about 5 m × 9 m and camera 0.6 m above the ground. (c, d) CARPM geometry and setup.

step, $p(\mathbf{x}_{t-1}|\mathbf{y}_{1:t-1})$, the motion-model, $p(\mathbf{x}_t|\mathbf{x}_{t-1})$ and the observation model, $p(\mathbf{y}_t|\mathbf{x}_t)$.

We use particle filtering methods in which the probability distribution of an unknown state is represented by a set of M weighted particles $\{\mathbf{x}_t^{(i)}, w_t^{(i)}\}_{i=1}^M$ [6]:

$$p(\mathbf{x}_t|\mathbf{y}_{1:t}) \approx \sum_{i=1}^M w_t^{(i)} \delta(\mathbf{x}_t - \mathbf{x}_t^{(i)}) \quad (7)$$

where $\delta(\cdot)$ is the Dirac delta function. Based on the discrete approximation of $p(\mathbf{x}_t|\mathbf{y}_{1:t})$, one can compute different estimates of the best state at time t . We use the Monte Carlo approximation of the expectation:

$$\hat{\mathbf{x}} \doteq \frac{1}{M} \sum_{i=1}^M w_t^{(i)} \mathbf{x}_t^{(i)} \approx \mathbb{E}(\mathbf{x}_t|\mathbf{y}_{1:t}), \quad (8)$$

or the maximum likelihood estimate:

$$\hat{\mathbf{x}}_{\text{ML}} \doteq \arg \max_{\mathbf{x}_t} \sum_{i=1}^M w_t^{(i)} \delta(\mathbf{x}_t - \mathbf{x}_t^{(i)}). \quad (9)$$

The tracking algorithm is composed by four steps:

1. *Prediction* – computes an approximation of $p(\mathbf{x}_t|\mathbf{y}_{1:t-1})$, by moving each particle according to the object's motion model.
2. *Observation* – computes the likelihood of each particle, based on image data.
3. *Update* – updates each particle's weight i using its likelihood $p(\mathbf{y}_t|\mathbf{x}_t^{(i)})$, by the means of $w_t^{(i)} \propto w_{t-1}^{(i)} p(\mathbf{y}_t|\mathbf{x}_t^{(i)})$.
4. *Resampling* – replicates the particles with a high weight and discards the ones with a low weight.

For this purpose, we need to model in a probabilistic way both the motion dynamics, $p(\mathbf{x}_t|\mathbf{x}_{t-1})$, and the computation of each particle's likelihood $p(\mathbf{y}_t|\mathbf{x}_t^{(i)})$ (steps 1 and 2). We discuss the model for the motion dynamics in the rest of this section, while we describe the observation model in Section 4.

3.1. Object motion dynamics

In order to accommodate to any real object motion, we use a Gaussian distribution, giving no privilege to any direction of motion:

$$p(\mathbf{X}_t|\mathbf{X}_{t-1}) = \mathcal{N}(\mathbf{X}_t|\mathbf{X}_{t-1}, \Lambda) \quad (10)$$

where $\mathcal{N}(\cdot|\mu, \Sigma)$ stands for a Gaussian distribution with mean μ and covariance Σ , and Λ stands for the diagonal matrix with the variances for random walk models on the components of the object state model. This approach has been widely used (e.g. [28,29]).

In this work we consider two kinds of objects: (i) spherical and (ii) polyhedral. For the first case, the state vector consists of the 3D Cartesian position and linear velocities of the ball, $\mathbf{X}_t = [x \ y \ z \ \dot{x} \ \dot{y} \ \dot{z}]^T$. The motion is modeled by a constant velocity model,

i.e., the motion equations correspond to a uniform acceleration during one time sample:

$$\mathbf{X}_{t+1} = \begin{bmatrix} \mathbf{I} & (\Delta t)\mathbf{I} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \mathbf{X}_t + \begin{bmatrix} \left(\frac{\Delta t^2}{2}\right)\mathbf{I} \\ (\Delta t)\mathbf{I} \end{bmatrix} a_t \quad (11)$$

where \mathbf{I} is the 3×3 identity matrix, $\Delta t = 1$, and a_t is a 3×1 white zero mean random vector corresponding to an acceleration disturbance. The covariance matrix of the random acceleration vector is usually set experimentally as $\text{cov}(a_t) = \sigma^2 \mathbf{I}$.

For the polyhedral object, the state vector is $\mathbf{X}_t = [\mathbf{p}_t; \mathbf{q}_t]$ where $\mathbf{p}_t = [x \ y \ z]^T$ denotes the position of the mass-center of the object and $\mathbf{q}_t = [q_w \ q_x \ q_y \ q_z]^T$ is a quaternion representing the object orientation. To model the dynamics, in this case we use a constant pose model, $\mathbf{p}_{t+1} = \mathbf{p}_t + \eta_p$ and $\mathbf{q}_{t+1} = \mathbf{q}_t * \eta_q$, where $*$ stands for quaternion product, η_p is Gaussian noise and η_q is a quaternion generated by sampling Euler angles from a Gaussian distribution.

Since the coordinates in the model are real-world coordinates, the motion model for a tracked object can be chosen in a principled way, both by using realistic models (constant velocity, constant acceleration, etc.) and by defining the covariance of the noise terms in intuitive metric units. The fact of using a constant velocity model is not limiting for cases where the objects can undergo sudden direction changes, e.g., in a RoboCup scenario. Using an adequately chosen acceleration noise, we can cope with arbitrary accelerations.

4. Observation model

In this section we describe the observation model, as expressed by $p(\mathbf{y}_t|\mathbf{x}_t)$ in (6). We propose a methodology that associates likelihood values to each of the samples in the particle filter using the target's 3D shape and color models. Recall that each particle represents an hypothesis of the target's position and pose. The likelihood function will assign high weights to a particle if the image information is coherent with its 3D position and pose hypothesis, and will be low otherwise. For instance, when tracking a red ball, an hypothesis whose projection is within a certain image region will have low likelihood if few red pixels are in that region. It will also be low if many pixels outside the region are equally red.

Fig. 2 illustrates this process. Boundaries of the multiple 3D pose hypothesis generated by the particle filter are projected onto the image plane with the sensor's projection function (we will explain latter that only some points inside and outside the boundary are projected, not the whole boundary). For visualization purposes, only a limited number of samples is shown – in practice a much higher number of samples is generated.

The likelihood function takes into account the 3D shape and the color properties of the targets. The rationale is to measure the similarity/dissimilarity between the colors of the target and the ones in the image regions corresponding to the boundaries of each hypothesis in the particle filter. High likelihood values should

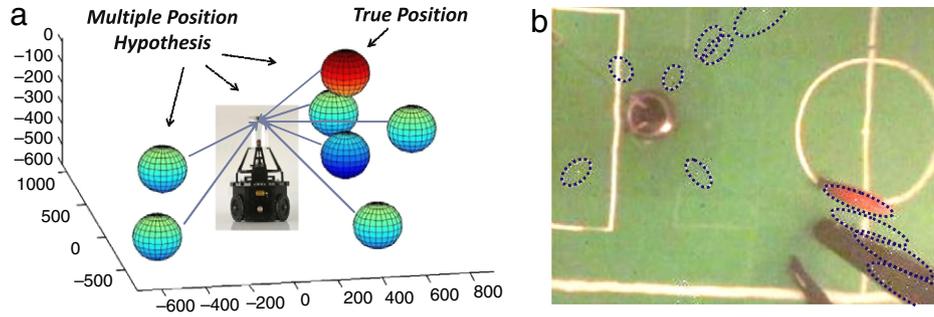


Fig. 2. The observation process projects the boundaries of the multiple 3D target hypothesis (a) into the image plane (b) and assigns likelihood values to each of them, as a function image color distributions inside and outside their boundaries.

only be assigned to particles if many pixels inside the region, and few outside, match the target's color model. High likelihood particles will survive and be replicated in the next time step while low likelihood particles will be rejected. An appropriately designed likelihood function is thus a crucial aspect in the convergence properties of the particle filter. It should be selective enough to quickly converge to a good approximation of the posterior (6). However it must not be too selective, otherwise particles close to the solution (but not close enough) will receive very low weights and will be rejected. Small mismatches in position and posture should nevertheless be accepted, despite having lower likelihood values than the optimal solution. We take these criteria into account in the definition of the likelihood function.

In practice we do not project the whole boundary contour in the image plane since boundary parameterizations may be hard to compute in arbitrary nonlinear sensors. Instead we define sets of 3D points, corresponding to specific regions of the target's surface or around it. These points are computed as a function of the object 3D shape and the particle's position and orientation with respect to the imaging system. These sets of 3D points, are projected onto the image plane, generating corresponding sets of points in 2D image coordinates. Then we build a color histogram for each set of points and compute the likelihood of the particle as a function of the similarities between pairs of color histograms. The only image information required in such process is the color of a few hundreds of pixels per particle. On the contrary to other existing methods, it only uses pixel values thus avoiding any other image processing operations, or rendering the full object model in the image plane. Also, it facilitates the utilization of nonlinear projection models, since only the projection of isolated points is required. In the consecutive sections we explain in depth each of the steps leading to the computation of particle's likelihood.

4.1. 3D points generation and their projection onto the image

From one 3D pose hypothesis for the tracked object we determine sets of 2D points that lie on the image, around the object edges and silhouette. The idea is that the color and luminance differences around the object edges are indicators of the likelihood of the hypothetical pose.

We consider two different object shapes: spheres and convex polyhedral objects. However the model can be easily extended to general polyhedra by exploiting the current knowledge in computer graphics [13].

For each state hypothesis $\mathbf{X}_t^{(i)}$, and given the particular 3D geometric object model, a few sets of 3D points are generated $\mathbf{U}_t^{(ij)}(\alpha)$, where t is the time step, i is the particle index, j indexes the points in the objects contour vicinity, and α represents a part of the object, usually the inside or outside of a surface. Then we convert the sets of 3D points $\mathbf{U}_t^{(ij)}(\alpha)$ to the corresponding sets of 2D points $\mathbf{u}_t^{(ij)}(\alpha)$ from which the likelihood measurements will be collected.

In the case of polyhedra, we use a 3D model that consists of a collection of faces and edges. To each pair (*face*, *edge*) we associate a set of 3D points that lie on the specific face, near the edge. To each edge we associate a set of points that lie on the corresponding edge of an expanded polyhedron (see Fig. 3b). The 3D points of this model are used to define the areas of the image where the color is sampled in order to build color histograms (see Fig. 3c and d). This is done by roto-translating the model and then projecting the 3D points onto the image.

For spheres, we define two sets of 3D points that when projected onto the image fall on the internal and external boundary of the sphere's silhouette. These 3D points lie on the intersection between the plane orthogonal to the line connecting the projection center to the center of the sphere and two spherical surfaces, one with a radius smaller than that of the tracked sphere, the other with a radius greater than that (see Fig. 3a).

The projection of the 3D points onto the image is performed using the appropriate projection model, as detailed in Section 2:

$$\mathbf{u}_t^{(ij)}(\alpha) = \mathcal{P}(\mathbf{U}_t^{(ij)}(\alpha)). \quad (12)$$

Fig. 4 shows examples of the obtained image points for a particle of the exact posture.

4.2. Color based likelihood measurement

The 2D points coordinates generated by the previous process are sampled in the current image, and their photometric information is used to obtain each particle's likelihood $w_t^{(i)}$. This approximates the state *a posteriori* probability density function, represented by the set of weighted particles $\mathbf{X}_t^{(i)}$, $w_t^{(i)}$. For color modeling we use independent normalized histograms in the HSI color space, which decouples intensity from color. We denote the B -bin reference (object) histogram model in channel $c \in \{H, S, I\}$ by $\mathbf{h}_{ref}^c = (h_{1,ref}^c, \dots, h_{B,ref}^c)$. An estimate for the histogram color model, denoted by $\mathbf{h}_x^c = (h_{1,x}^c, \dots, h_{B,x}^c)$, can be obtained as

$$h_{i,x}^c = \beta \sum_{\mathbf{u} \in \mathcal{U}} \delta_a(b_{\mathbf{u}}^c), \quad i = 1, \dots, B. \quad (13)$$

\mathcal{U} is the region where the histogram is computed; $b_{\mathbf{u}}^c \in \{1, \dots, B\}$ denotes the histogram bin index associated with the intensity at pixel location \mathbf{u} in channel c ; δ_a is a Kronecker delta function at a ; and β is a normalization constant such that $\sum_{i=1}^B h_{i,x}^c = 1$.

We define $\mathbf{h}^{\text{model}}$, \mathbf{h}^{in} and \mathbf{h}^{out} as a reference (object) histogram, the inner boundary points and the outer boundary points histogram, respectively. We define \mathbf{h}_{sideA}^i and \mathbf{h}_{sideB}^i as the histograms of each of the two sides of the i th non-silhouette edge (see Fig. 3d). To measure the difference between histograms we use the Bhattacharyya similarity as in [30,31]:

$$\mathcal{S}(\mathbf{h}_1, \mathbf{h}_2) = \sum_{i=1}^B \sqrt{h_{i,1} h_{i,2}}. \quad (14)$$

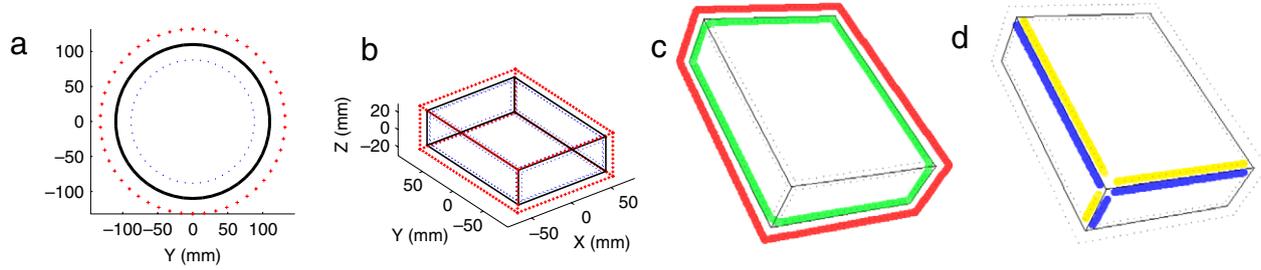


Fig. 3. Sphere and polyhedron 3D model. (a) Location of the internal (blue dots) and external (red crosses) 3D points sets for the sphere, with respect to the great circle of the sphere (black line). (b) Location of the 3D points of the polyhedron model: the points lying on the faces of the object are represented as blue dots, while the points lying on the edges of the expanded object are represented as red crosses. (c) One particular projection of the polyhedron. The areas sampled in the image to build the internal and external histograms are highlighted. (d) The same projection of the polyhedron with the pairs of areas associated with non-silhouette edges highlighted. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

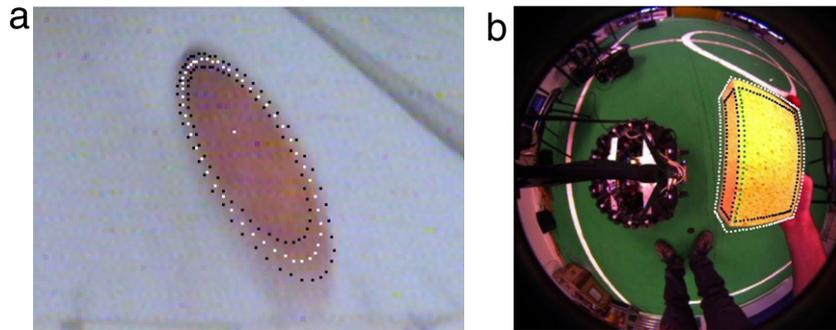


Fig. 4. Examples of image points used to build the color histograms used for particle likelihood computation. Obtained from the projection of the 3D points shown in Fig. 3. (a) For spheres. (b) For cuboids.

We define:

$$\begin{aligned} \mathcal{S}_0 &= \mathcal{S}(\mathbf{h}^{\text{model}}, \mathbf{h}^{\text{in}}), & \mathcal{S}_1 &= \mathcal{S}(\mathbf{h}^{\text{out}}, \mathbf{h}^{\text{in}}), \\ \mathcal{S}_2 &= \frac{\sum_{i=0}^n \mathcal{S}(\mathbf{h}_i^{\text{sideA}}, \mathbf{h}_i^{\text{sideB}})}{n} \end{aligned} \quad (15)$$

as the object-to-model, object-to-background and mean-side-to-side (non-silhouette edges) similarities, respectively. Finally, the resulting distance is:

$$\mathcal{D} = \left(1 - \frac{\mathcal{S}_0 + \kappa_1(1 - \mathcal{S}_1) + \kappa_2(1 - \mathcal{S}_2)}{\kappa_1 + \kappa_2 + 1} \right) - \gamma \quad (16)$$

where γ is a coefficient that modulates the distance based on the number of projected 3D points that fall onto the image, $\gamma = \ln\left(\frac{\text{used points ratio}}{\epsilon}\right)$.

The rationale for this definition of \mathcal{D} is that the distance metric should be high when candidate color histograms are different from the reference histogram and similar to the background. It should also be high when there is little or no difference in color on the sides of non-silhouette edges. Parameters κ_1 and κ_2 allow to balance the influence of the different contributions, up to the extent of ignoring them, by setting such parameters to 0. They were set to 1.5 and 0.6 respectively, for the case of the polyhedron; for tracking the sphere they were set to 1.5 and 0 in the first and second experiment and to 0 and 0 in the third one (thus ignoring the reference color model). The term γ is useful when tracking an object whose projection lies only partly on the image. The data likelihood function \mathcal{L} is modeled as a Laplacian distribution over the distance metric: $p(\mathbf{y}_t | \mathbf{X}_t^{(i)}) \propto e^{-\frac{|\mathcal{D}|}{\epsilon}}$. In our experiments we set $\epsilon = 1/30$.

5. Experimental results

This section presents an evaluation of the proposed methods. Firstly we present results taken with omnidirectional cameras:

the tracking of a ball performed with a catadioptric setup and the tracking of a cuboid in a dioptric setup. Secondly we present results with conventional cameras and perform a comparison between our method and a competing alternative based on 2D tracking followed by 3D reconstruction. In this set of experiments we show results with real and artificial images, both with ground truth. Finally we show results taken from a mobile platform equipped with a dioptric vision system, tracking a ball of unknown color.

The first and the last sets of experiments are taken in very unconstrained scenarios and their evaluation is not supported by ground truth values of the target's trajectories. Evaluation of the results is made only on a quantitative basis by empirical observations of the system's performance. On the contrary, the second set of experiments is supported by ground truth values and results are evaluated quantitatively with the mean squared tracking error. In the experiments with real images, ground truth is obtained manually by measuring the 3D coordinates of some points along the trajectory. In the experiments with artificial images, target's true 3D coordinates are provided by the simulation engine at each time step.

5.1. Omnidirectional cameras

Here we illustrate qualitatively the performance of the methods with omnidirectional images with large nonlinear distortions, blur and noise. The sequences were collected from setups usually employed in RoboCup Middle Size league scenarios.

5.1.1. Tracking a ball in 3D with one catadioptric omnidirectional camera

In this experiment we tracked a ball hitting an obstacle on the ground and subsequently performing a series of parabolic movements. This image sequence was acquired with a Marlin firewire camera (see Fig. 1b) with a frame rate of 25 fps and a resolution

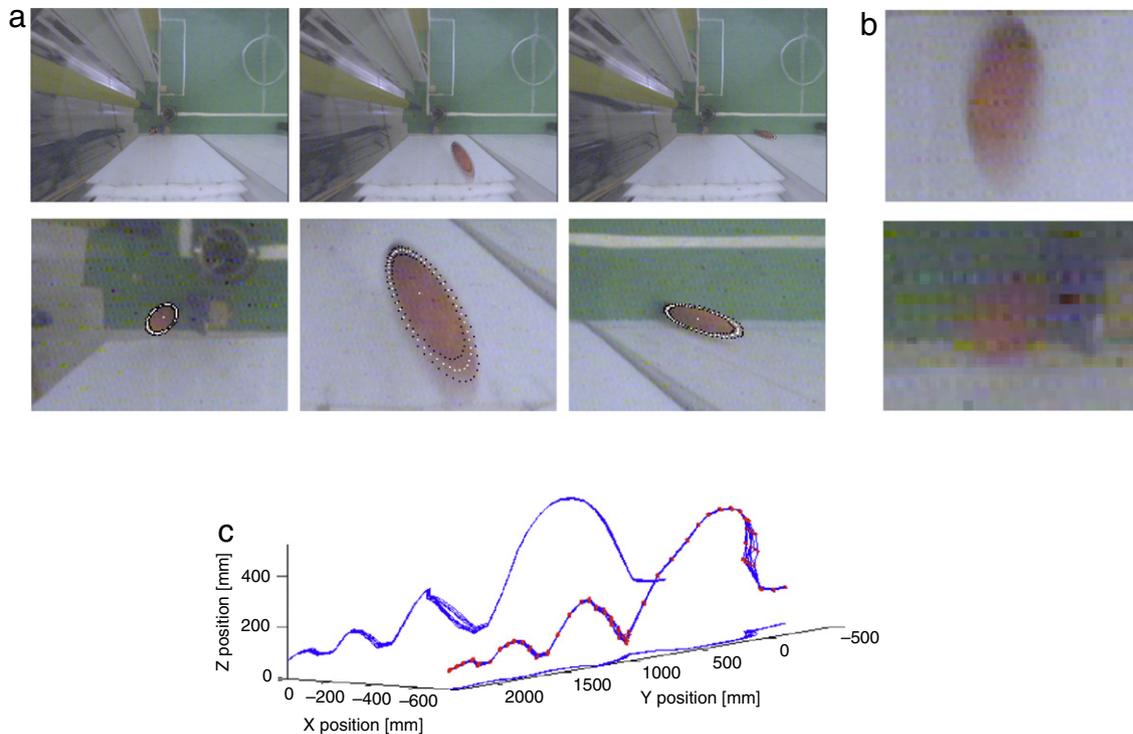


Fig. 5. Ball tracking in the catadioptric setup. (a) Three frames of the sequence (top row), the corresponding close-ups (bottom row). In both cases the projection of a ball lying in the estimated position is marked in white, while the pixels used to build the inner and outer color histograms are marked in black. (b) Close-ups of the ball showing motion blur and image sensor noise. (c) Plot of the tracked paths resulting from 10 runs of the algorithm performed on the same image sequence. The 10 blue lines with red points represents the 10 3D estimated trajectories of the ball, the blue lines are the projection of these trajectories on the ground and lateral plane. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

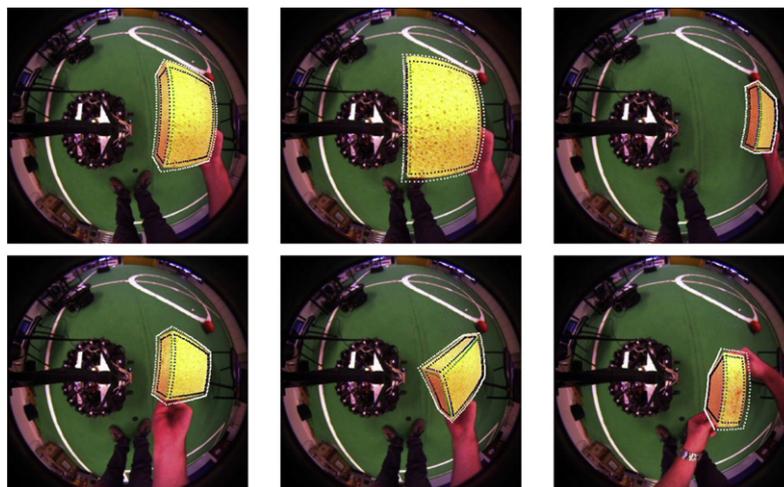


Fig. 6. Cuboid tracking in the dioptic setup. Six frames of the sequence with the pixels used to build color histograms highlighted.

of 640×480 pixels. The tracker used 10000 particles, the Unified Projection Model described in Section 2 and the constant velocity motion model described in Section 3.1. Processing was done off-line, with Matlab. Processing one frame took 11 s on a 2 GHz Pentium 4 CPU, i.e. 1.1 ms per particle. In this case the tracker was provided with the color model for the ball. The projection of the ball on the image plane changes dramatically in size during the image sequence (see Fig. 5a), due to the nature of the catadioptric system used. The images are affected by both motion blur and heavy sensor noise (see Fig. 5b). We repeated the tracking 10 times on the same image sequence to assess how the stochastic component of the tracker influences the precision of the 3D estimated paths (see Fig. 5c), with satisfactory results.

5.1.2. Tracking a cuboid in 3D with one dioptic omnidirectional camera

In this experiment we tracked a yellow cuboid in a sequence of 600 frames, acquired using the dioptic omnidirectional setup, comprising a Marlin firewire camera and a Sunex DSL215 lens (see Fig. 1d). The resolution was 640×480 pixels. In the tracker we used 5000 particles, the Constant Angular Resolution Projection Model described in Section 2 and the constant position motion model described in Section 3.1. Processing was done off-line, with Matlab. Processing one image took 33 s on a 2 GHz Pentium 4 CPU, i.e. 6.6 ms per particle. The tracker managed to follow the cuboid along rotations and translations that greatly affect its projection onto the image, see Fig. 6.



Fig. 7. Non-uniformly colored ball tracking on cluttered background. Four frames of the sequence with the estimated contour highlighted in green (top row). Close-ups on the ball showing the behaviour of the tracker in challenging conditions: partial occlusion, heavy background clutter, motion blur, projection of the ball lying partly outside the image (middle row). One close-up showing the behaviour of the tracker when only a small part of the projection of the ball lies on the image. Two frames showing the different illumination conditions at the beginning and at the end of the ramp. A plot of the intensity histogram built on the area occupied by the ball at the beginning of the ramp, in red, and at the end of the ramp, in blue (bottom row). (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

5.2. Comparative study: perspective cameras with ground truth

Here we evaluate our algorithm quantitatively with ground truth data and compare its performance against an alternative approach: a 2D tracker followed by a 3D reconstruction algorithm. The 2D tracker is based on a MMDA (Multiple Model Data Association) methodology [32] and outputs a collection of points uniformly sampled at the target's contour. In each image, the tracked contour is expanded in the 3D cone generated by the camera optical center and the sphere's occluding contour. Then, knowing the sphere's radius, simple trigonometric calculations allow the reconstruction of its 3D position.

5.2.1. Tracking a non-uniformly colored ball in 3D on cluttered background

In this experiment we tracked a non-uniformly colored ball rolling down a ramp. The path of the ball consists of four different parts: in the first part the ball is on top of the ramp and partially occluded by the hand of an operator, in the second part the ball rolls down the ramp, in the third part the ball performs a free fall, while in the fourth it is caught by an operator and is gradually moved out of the image. The image sequence is affected by heavy background clutter, occlusions, motion blur and changes in illumination (see Fig. 7). It was acquired with a Canon XM1 camera with a frame rate of 25 fps and a resolution of 640×480 pixels. The tracker used 10000 particles, a version of the Perspective Projection Model that takes into account radial and tangential distortion, and the constant velocity motion model described in Section 3.1. Processing was done off-line, with Matlab. Processing one image took 21 s on a 2 GHz Pentium 4 CPU, i.e. 2.1 ms per particle. The tracker was provided with the color model for the ball.

The proposed method successfully tracked the ball along the whole sequence (see Fig. 8a). We compared the performance of the proposed tracker with the MMDA-based method, finding that

Table 1

Precision of the 3D tracking methods (PF – Particle Filter, MMDA – Multiple Model Data Association) during the linear part of the ramp descent experiment. Precision is measured by the root mean squared error (RMSE) between the measured trajectory and the best fit linear trajectory. Notice that this error metric does not take into account constant bias terms (systematic errors).

Method	Precision (RMSE)
PF	0.035 m
MMDA	0.017 m

they produce similar results (see Fig. 8b). This comparison, though, could only be performed on the part of the sequence in which the ball is fully visible, as the MMDA tracker is not robust to occlusions.

The proposed method is able to cope with occlusions, although the 3D estimate is less accurate: the gradual occlusion during the last part of the path induces an increasing localization error, but that is difficult to quantify as the motion of the ball was constrained by an operator catching it.

To evaluate the precision of both methods during the linear part of the ramp descent, we have computed the root-mean-squared error of the estimated position with respect to a “ground truth trajectory” obtained by fitting a line to the observed 3D points. Results are presented in Table 1. Furthermore, we compared the performance of the two methods along the first and second half of the linear path, assessing that the increasing motion blur does not affect it. The RMS error for the PF approach is smaller during the second half of the linear path (0.039 m and 0.032 m, for the first and second halves respectively), as it is for the MMDA method (0.019 m and 0.016 m, respectively). Eventually, we measured the precision of the PF approach during the frames in the beginning of the sequence while the ball is partially occluded and still, obtaining a value of 0.033 m.

We notice that both methods have low tracking errors. The PF filter is less precise than the MMDA method in this particular

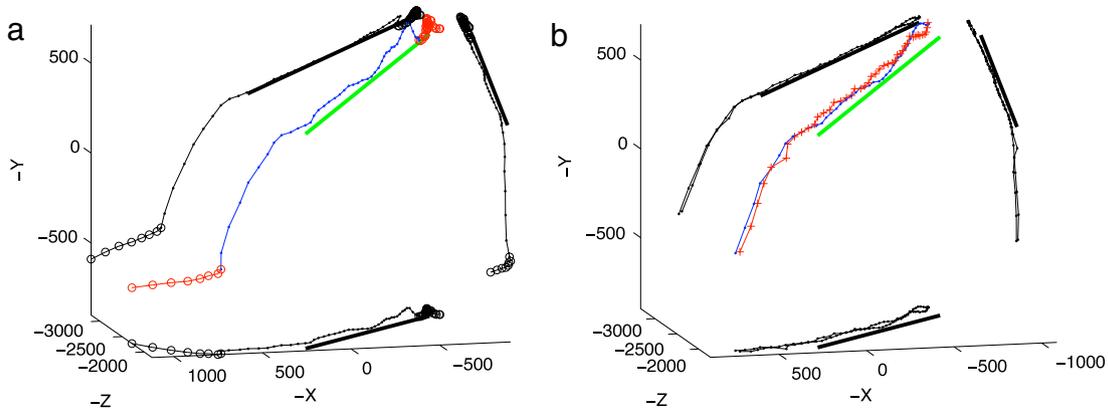


Fig. 8. (a) 3D plot of the trajectory estimated with PF along the whole image sequence. The parts of the trajectory during which the ball is occluded are plotted with red circled lines, while the part in which the ball is fully visible is plotted as a blue dotted line. Ground truth along the ramp is represented with a green line, the projection of each line on the ground, lateral and vertical plan is shown in black. (b) 3D plot of the trajectories estimated with PF (blue dotted line) and the MMDA-based tracker (red crossed line). Because the MMDA-based method can not deal with occlusion, both the beginning and the end of the sequence (where occlusions are significant and the MMDA tracker fails) are not shown in this plot. Ground truth is plotted in green, projections on the plans are plotted in black. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

Table 2

Accuracy of the 3D tracking methods (PF – Particle Filter, MMDA – Multiple Model Data Association) for the virtual ball tracking experiment. Accuracy is measured by the root mean squared error (RMSE) between the measured trajectory and the true synthesized trajectory. Both systematic errors (bias) and random errors (variance) are taken into account on the error metric.

Method	Accuracy (RMSE)
PF	0.0192 m
MMDA	0.0712 m

experiment, this is due to an arbitrary choice of the tuning parameters. By appropriately tuning the acceleration covariance noise it would be possible to change the precision of the methods. What we would like to stress in this experiment is that the PF method is able to cope better with significant levels of occlusion whereas the MMDA method fails.

5.2.2. Tracking a non-uniformly colored ball in 3D in simulated images with ground truth

In this experiment we tracked a non-uniformly colored ball along a spiraling path. The images were generated with a virtual reality software,¹ so the intrinsic parameters of the camera and the 3D position of the ball at each given time are known exactly. The image sequence was generated with a resolution of 576×380 pixels. The tracker used 10000 particles, the Perspective Projection Model, and the constant velocity motion model. Processing was done off-line, with C++. Processing one image took 0.26 s on a 2 GHz Pentium 4 CPU, i.e. 0.026 ms per particle. The tracker was provided with the color model for the ball. The results from simulation experiments can be seen in Fig. 9.

The PF and MMDA methods were compared in terms of root-mean-squared error to the true trajectory. Whereas in the previous experiment we were able to evaluate only the precision of the methods (repeatability), in this experiment we are able to quantify accuracy (truthfulness).

In comparison to the MMDA-based method, we notice a better match with the ground truth, despite a slightly larger variance. Quantitatively, the average error is much smaller with the proposed tracking method. Table 2 presents the obtained accuracy values.

The MMDA technique was designed to track approximately circular shapes in the image. Due to a large focal distance of

the camera used in this experience, there is a high perspective distortion and the 3D sphere projects to the image plane into an ellipse whose eccentricity varies significantly with position. This produces a space variant bias in the estimated 3D position that was not so obvious in the previous experiment. On the contrary, our method does not suffer from this effect because it uses the true 3D space invariant shape of the object instead of its 2D space variant projection. It is, therefore, more easily applicable to different kinds of sensors.

5.3. Arbitrary color ball tracking with a moving robot

To use the tracker on moving robots we need to distinguish between two reference frames for the 3D coordinates of a point: the world reference frame $\{W\}$, which is inertial, and the robot reference frame $\{R\}$, which is not inertial as the robot undergoes accelerations during its motion. The state of a tracked object is expressed in terms of the world reference frame, so that a motion model based on the laws of physics can be expressed in the simplest possible way (Fig. 10).

In order to project a 3D point onto the image plane, its relative position with respect to the imaging system of the robot must be known. This means that coordinates expressed in the world reference frame ${}^W M = [{}^W X, {}^W Y, {}^W Z, 1]$ must be transformed to the robot reference frame ${}^R M = [{}^R X, {}^R Y, {}^R Z, 1]$ by means of the transformation matrix ${}^R T_W$, which comprises the rotation matrix ${}^R R_W$ and the translation vector ${}^R t_W$.

$${}^R T_W = \begin{bmatrix} {}^R R_W & {}^R t_W \\ 0 & 1 \end{bmatrix}. \quad (17)$$

The transformation matrix is computed at every time step based on the difference between the initial pose of the robot and the current one. The 2D coordinates of the point m , corresponding to the point ${}^W M$, are thus computed as:

$$m = \mathcal{P}({}^R M) = \mathcal{P}({}^R T_W \cdot {}^W M). \quad (18)$$

In other words, we model the dynamics of the tracked object in the world reference frame, while the observer position is taken into account in the observation model.

The tracker was implemented in the soccer robot team ISocRob [33] software architecture in order to demonstrate a real-time application of the method with robots playing in the RoboCup Middle Size League. In this experiment the omnidirectional robot tracked a moving ball, moved around by a human. The robot is constantly moving towards the ball without catching it.

¹ The authors would like to thank Simon Day for granting them the possibility to use one of his images as wallpaper in this experiment.

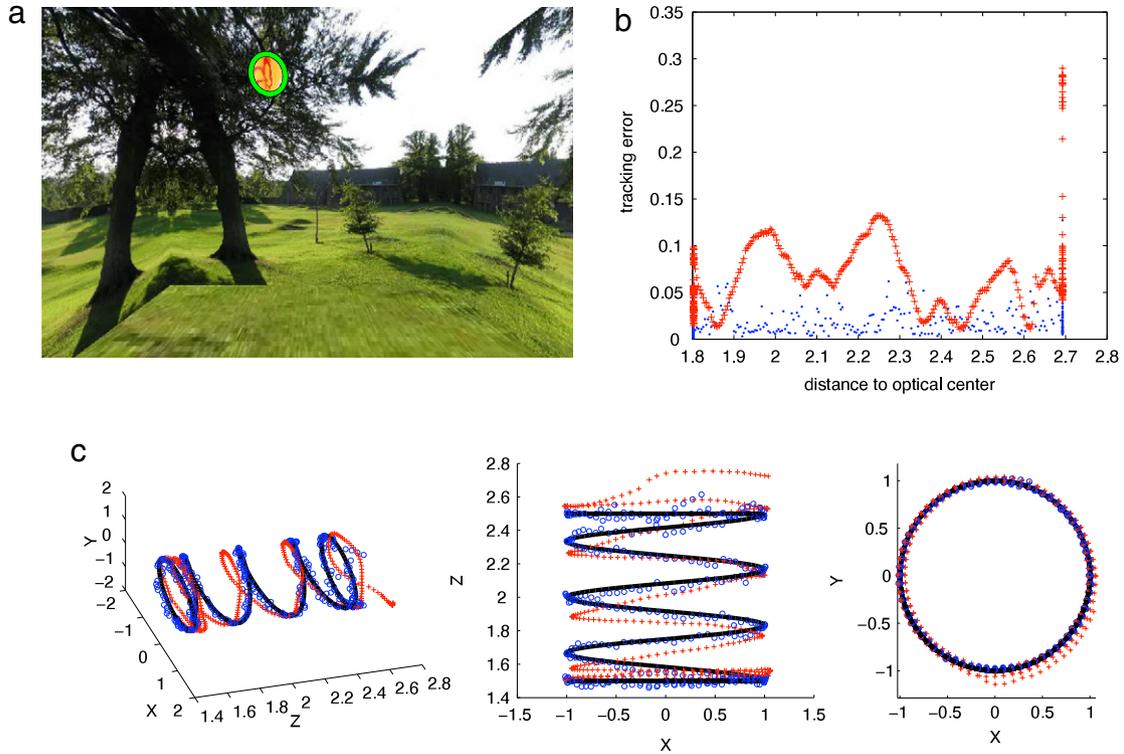


Fig. 9. Non-uniformly colored ball tracking with ground truth. (a) One frame of the sequence with the estimated contour highlighted in green. (b) Plot of the tracking error as a function of distance to the projection center for PF (blue dots) and the MMDA-based approach (red crosses). (c) Plot of the tracked path of the ball estimated with PF (blue dots) and the MMDA-based approach (red crosses). Each graph represents a different view of the reconstructed 3D trajectory. The left plot is a general perspective view, the middle plot is a projection in X, Z plane (top view) and the right plot is a projection in the X, Y plane (front view). The ground truth trajectory is represented by the black line. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

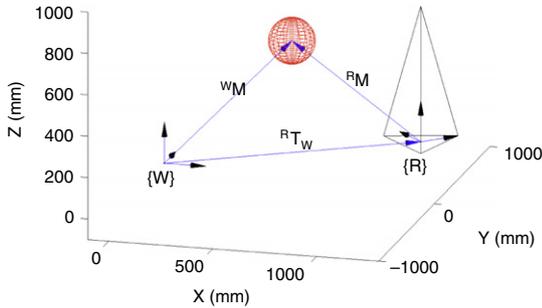


Fig. 10. Position of one sphere in the world and robot reference frames.

In this implementation the observation model discards the object-to-model mismatch, described in Eq. (15), and relies strictly on the object-to-background dissimilarity. Therefore, we carried out the experience with an ordinary soccer ball, mainly white colored, as one can see in Fig. 11a, but other colored balls, e.g., orange, could be used.

Images on the robot were acquired with the dioptic omnidirectional camera used in Section 5.1, at 10 fps. Odometry motion control measurements were obtained at 25 fps and we used only 600 particles in the tracker. Processing was done on-line, at 10 fps. The results are visible in Fig. 11b.² In the figure, the robot posture and ball position are plotted in a frame located on the center of the soccer field (*field frame*). The ball coordinates were transformed from the *robot frame* to the *field frame* and the robot self-localization postures (obtained using Monte Carlo Localization [34] – MCL) were

expressed in the *field frame*, to plot the ball at each update step. The ball trajectory is not smooth due to the robot self-localization errors. Ball shifts in the plot correspond to equal shifts in the robot posture, when MCL corrects the robot posture, since the transformation of the ball position from the *robot frame* to the *field frame* is based on the current estimate of the robot posture by MCL.

6. Conclusions

We presented a 3D model based tracking system, based on a particle filter framework. The method requires the knowledge of the 3D shape of the target and the imaging sensor calibration. We stress that the system is particularly suited for omnidirectional vision systems, as it only requires the projection of isolated points arising from likely posture hypotheses for the target. In practice the method can be used with any projection model given the availability of ways to compute the projection of 3D points to sensor coordinates. We have shown applications with different kinds of omnidirectional vision sensors (dioptic and catadioptric), with different objects (spheres and polyhedra) and with moving observers. The proposed methods were evaluated in challenging real and synthetic sequences with clutter, partial occlusions, image blur and non-uniformly colored objects. The observation model assumes distinct target and background color distributions and benefits, but does not depend critically, on their knowledge. We have shown a case where the target color model is absent and tracking relies only on the contrast between target and background color. Anyway having known and distinctly colored objects, naturally improves the robustness of the tracker and makes it specially suited to color engineered environments as is the case of robotic competitions.

The proposed method was compared with a state-of-the-art 2D tracker followed by 3D pose estimation, in the problem of

² See also a video taken online – at 2 fps – from the robot camera during this same experiment, at <http://www.youtube.com/watch?v=szReXrJueYQ&feature=channel>.

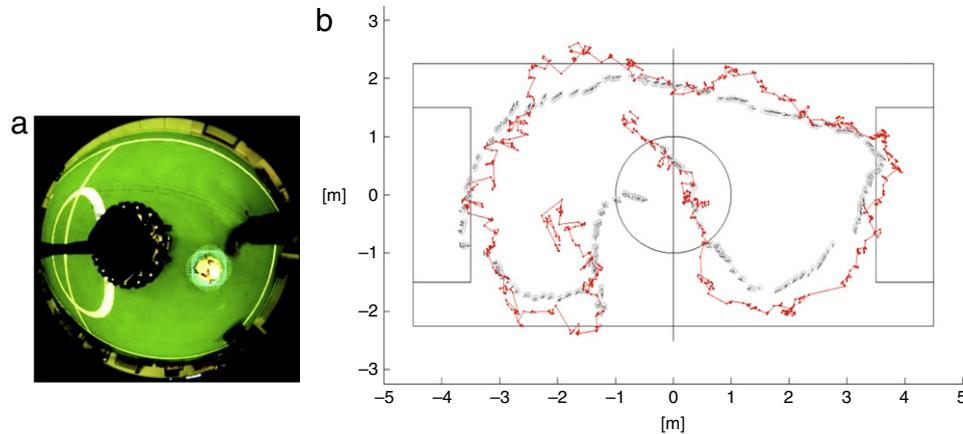


Fig. 11. Tracking a white ball with a moving robot. (a) Detection of a white ball based on the object-to-background dissimilarity, where the light green crosses mark the pixels used to build the background histogram. (b) Plot of the paths of robot and ball. The robot starts in the middle circle facing opposite to the ball. The gray circles represent the robot's pose while the red dots represent the ball localization, here in a 2D representation only. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

tracking spheres. The two methods have similar results whenever the 2D tracker reliably computes the image contour. However, the proposed methodology has proven more robust to partial occlusion. Additionally, our method is advantageous in terms of the customization required to cope with objects of different shapes. In our case we just have to recompute points in the inner and outer object boundaries, compute their visibility and project them to the sensor plane (a large body of computer graphics algorithms and hardware acceleration techniques are available for such purposes). Instead, contour based methods require 3D pose estimation from the 2D silhouette which often involves complex computations and is frequently an ill-posed problem.

In future work we will advance in two directions. First we will consider the use of multiple models for target motion. Despite a constant velocity model, as the one employed in this work, can cope with arbitrary accelerations by increasing the model's acceleration noise, this has the side effect of increasing also the variance of the estimated trajectories. With multiple motion models we expect to have less noisy position estimates and still be able to cope with sudden target accelerations. Secondly we will investigate the possibility of applying information fusion methods to perform the collaborative tracking of objects in distributed robotics systems. Something we did not address in this paper was the consideration of localization errors in the robot/sensor position information. By appropriately combining the information of both robots' and targets' location and uncertainty, we aim at attaining a globally coherent and precise localization of robot teams.

Acknowledgements

We would like to thank Dr. Luis Montesano, Dr. Alessio Del Bue and Giovanni Saponaro for the fruitful discussions.

References

- [1] P. Lima, A. Bonarini, C. Machado, F. Marchese, F. Ribeiro, D. Sorrenti, Omnidirectional catadioptric vision for soccer robots, 2001.
- [2] J. Gaspar, N. Winters, J. Santos-Victor, Vision-based navigation and environmental representations with an omnidirectional camera, *IEEE Transactions on Robotics and Automation* 16 (6) (2000).
- [3] T. Boulton, X. Gao, R. Michaels, M. Eckmann, Omnidirectional visual surveillance, *Image and Vision Computing* 22 (7) (2004) 515–534.
- [4] Yasuo Kuniyoshi, Nobuyuki Kita, Sebastien Rougeaux, Takashi Suehiro, Active stereo vision system with foveated wide angle lenses, in: *ACCV'95: Invited Session Papers*, Springer-Verlag, London, UK, 1996, pp. 191–200.
- [5] V. Lepetit, P. Fua, Monocular model-based 3d tracking of rigid objects: A survey, *Foundations and Trends in Computer Graphics and Vision* 1 (1) (2005) 1–89.
- [6] A. Doucet, N. de Freitas, N. Gordon, in: *Gordon (Ed.), Sequential Monte Carlo Methods in Practice*, Springer Verlag, 2001.
- [7] M. Isard, A. Blake, Condensation: conditional density propagation for visual tracking, *International Journal of Computer Vision* 28 (1) (1998) 5–28.
- [8] Y. Wu, Robust visual tracking by integrating multiple cues based on co-inference learning, *International Journal of Computer Vision* 58 (1) (2004) 55–71.
- [9] Z. Khan, T. Balch, F. Dellaert, A rao – blackwellized particle filter for eigentracking, in: *Proc. of Int. Conf. on Computer Vision and Pattern Recognition, CVPR'04*, vol. 2, 2004, pp. 980–986.
- [10] Changjiang Yang, Ramani Duraiswami, Larry Davis, Fast multiple object tracking via a hierarchical particle filter, in: *Proc. of Int. Conf. on Computer Vision and Pattern Recognition, CVPR'05*, 2005, pp. 212–219.
- [11] Ya-Dong Wang, Jian-Kang Wu, Ashraf A. Kassim, Particle filter for visual tracking using multiple cameras, in: *Proc. of IAPR Conf. on Machine Vision Applications, MVA'05*, 2005, pp. 298–301.
- [12] K. Choo, D.J. Fleet, People tracking using hybrid monte carlo filtering, in: *Proc. Int. Conf. on Computer Vision, ICCV'01*, 2001 pp. 321–328.
- [13] G. Klein, D. Murray, Full-3d edge tracking with a particle filter, in: *Proc. of BMVC 2006*, Edinburgh, Scotland, 2006, pp. III:1119.
- [14] M. Pupilli, A. Calway, Real-time camera tracking using known 3d models and a particle filter, in: *Proc. of ICPR (1)*, Hong Kong, 2006, pp. 199–203.
- [15] Tatsuya Osawa, Xiaojun Wu, Kaoru Wakabayashi, Takayuki Yasuno, Human tracking by particle filtering using full 3d model of both target and environment, in: *ICPR (2)*, IEEE Computer Society, 2006, pp. 25–28.
- [16] L. Montesano, M. Lopes, A. Bernardino, J. Santos-Victor, Learning object affordances: From sensory motor maps to imitation, *IEEE Transactions on Robotics* 24 (1) (2008) 15–26.
- [17] M. Taiana, J. Gaspar, J. Nascimento, A. Bernardino, P. Lima, 3D tracking by catadioptric vision based on particle filters, in: *Proc. of the Robocup Symposium*, Atlanta, 2007.
- [18] M. Tajana, J. Gaspar, J. Nascimento, A. Bernardino, P. Lima, On the use of perspective catadioptric sensors for 3d model-based tracking with particle filters, in: *IEEE/RSJ Int. Conf. on Intelligent Robots and Systems*, 2007 pp. 2747–2752.
- [19] M. Taiana, J. Nascimento, J. Gaspar, A. Bernardino, Sample-based 3D tracking of colored objects: a flexible architecture, in: *Proc. of BMVC 2008*, Leeds, 2008.
- [20] M. Taiana, J. Santos, J. Gaspar, J. Nascimento, A. Bernardino, P. Lima, Color 3D model-based tracking with arbitrary projection models, in: *Workshop Proceedings of SIMPAR 2008: Omnidirectional Robot Vision*, Venice, Italy, 2008, pp. 206–217.
- [21] H. Bakstein, T. Pajdla, Panoramic mosaicing with a 180 deg field of view lens, in: *IEEE Workshop on Omnidirectional Vision*, 2002, pp. 60–67.
- [22] B. Micusik, T. Pajdla, Estimation of omnidirectional camera model from epipolar geometry, in: *IEEE CVPR*, 2003, pp. 485–490.
- [23] R. Hicks, R. Bajcsy, Catadioptric sensors that approximate wide-angle perspective projections, in: *IEEE CVPR*, 2000, pp. 545–551.
- [24] J. Gaspar, C. Deccó, Jun Okamoto Jr, J. Santos-Victor, Constant resolution omnidirectional cameras, in: *IEEE Workshop on Omni-directional Vision*, 2002 pp. 27–34.
- [25] G. Geyer, K. Daniilidis, A unifying theory for central panoramic systems and practical applications, in: *IEEE CVPR*, 2000, pp. 445–461.
- [26] S. Baker, S.K. Nayar, A theory of single-viewpoint catadioptric image formation, *International Journal of Computer Vision* 35 (2) (1999) 175–196.
- [27] R. Benosman, S. Kang, *Panoramic Vision*, Springer Verlag, 2001.
- [28] P. Barrera, J.M. Cañas, V. Matellán, Visual object tracking in 3d with color based particle filter, in: *Proc. of World Academy of Science, Eng. and Tech.*, vol. 4, 2005.
- [29] P. Pérez, J. Vermaak, A. Blake, Data fusion for visual tracking with particles, *Proceedings of the IEEE* 92 (3) (2004) 495–513.
- [30] D. Comaniciu, V. Ramesh, P. Meer, Real-time tracking of non-rigid objects using mean shift, in: *Proc. Conf. on Comp. Vision Pattern Rec.*, vol. 2, 2000, pp. 142–149.

- [31] K. Okuma, A. Taleghani, N. de Freitas, J.J. Little, D.G. Lowe, A boosted particle filter: multitarget detection and tracking, in: Proc. ECCV, 2004, pp. 28–39.
- [32] J. Nascimento, J.S. Marques, Robust shape tracking with multiple models in ultrasound images, *IEEE Transactions on Image Processing* 17 (3) (2008) 392–406.
- [33] H. Costelha, N. Ramos, J. Estilita, J. Santos, M. Taiana, J. Torres, T. Antunes, P. Lima, Isocrob 2007 team description paper, Technical Report, Instituto Superior Técnico, 2007.
- [34] J. Messias, J. Santos, J. Estilita, P. Lima, Monte Carlo localization based on gyrodometry and line-detection, in: ROBOTICA2008, 8th Conference on Mobile Robots and Competitions, 2008.



M. Taiana received his M.Sc. degree in Computer Engineering from Politecnico di Milano - Italy, in 2007. He is currently a Ph.D. student at the Computer Vision Laboratory (VisLab), which belongs to the Institute for Systems and Robotics (ISR) of Instituto Superior Técnico (IST) - Lisbon. His research interests are in Computer and Robot Vision and Robotics.



J. Santos received his M.Sc. degree in Electrotechnical and Computer Engineering from Instituto Superior Técnico - Portugal, in 2008. He is co-founder at selfTech - Engenharia de Sistemas e Robótica, a spin-off of Institute for Systems and Robotics (ISR), currently focused in deploying autonomous robotic solutions. His research interests are in Mobile Robot Localization and Computer Vision.



J. Gaspar, received his Ph.D. degree in Electrical and Computer Engineering from Instituto Superior Técnico (IST), Technical University of Lisbon - Portugal, in 2003. He is currently an Auxiliary Professor at the IST, and a Researcher at the Computer Vision Laboratory (VisLab), Institute for Systems and Robotics (ISR). His research interests are in Computer and Robot Vision, Robotics and Control.



J. Nascimento (M'06) received the EE degree from Instituto Superior de Engenharia de Lisboa, in 1995, and the M.Sc. and Ph.D. degrees from Instituto Superior Técnico (IST), Technical University of Lisbon, in 1998, and 2003, respectively. Presently, he is a postdoctoral researcher with the Institute for Systems and Robotics (ISR) at IST. His research interests include image processing, shape tracking, robust estimation, medical imaging, and video surveillance. Dr. Nascimento has co-authored over 50 publications in international journals and conference proceedings (many of which of the IEEE), has served on program committees of many international conferences, and has been a reviewer for several international journals.



A. Bernardino received the Ph.D. degree in Electrical and Computer Engineering in 2004 from Instituto Superior Técnico (IST). He is an Assistant Professor at IST and Researcher at the Institute for Systems and Robotics (ISR-Lisboa) in the Computer Vision Laboratory (VisLab). He participates in several national and international research projects in the fields of robotics, cognitive systems, computer vision and surveillance. He published several articles in international journals and conferences, and his main research interests focus on the application of computer vision, cognitive science and control theory to advanced robotic and surveillance systems.



P. Lima got his Ph.D. (1994) in Electrical Engineering at the Rensselaer Polytechnic Institute, Troy, NY, USA. Currently, he is an Associate Professor at Instituto Superior Técnico, Lisbon Technical University. He is also a member of the Institute for Systems and Robotics, a Portuguese research institute, where he is coordinator of the Intelligent Systems group and a member of the Scientific Board. Pedro Lima is a Trustee of the RoboCup Federation, and was the General Chair of RoboCup2004, held in Lisbon. He currently serves as President Portugal Robotics Society, for the 2009/10 period. He is also a founding member of

Portugal Robotics Society (2006), of the IEEE RAS Portugal Chapter (2005), and a senior member of the IEEE. He is the co-author of two books, regularly serves as member of international conferences program committees, and has coordinated national and international (ESA, EU) R&D projects. He has also been very active in the promotion of Science and Technology to the society, through the organization of Robotics events in Portugal, including the Portuguese Robotics Open since 2001. His research interests lie in the areas of discrete event models of robot plans, reinforcement learning and planning under uncertainty, with applications to multi-robot systems.