

# Angle-based Homing from a Reference Image Set using the 1D Trifocal Tensor

M. Aranda, G. López-Nicolás and C. Sagüés

Instituto de Investigación en Ingeniería de Aragón, Universidad de Zaragoza, Spain

Email: {marandac, gonlopez, csagues}@unizar.es

**Abstract**— This paper presents a visual homing method for a robot moving on the ground plane. The approach employs a set of omnidirectional images acquired previously at different locations (including the goal position) in the environment, and the current image taken by the robot. We present as contribution a method to obtain the relative angles between all these locations, using the computation of the 1D trifocal tensor between views and an indirect angle estimation procedure. The tensor is particularly well suited for planar motion and provides important robustness properties to our technique. Another contribution of our paper is a new control law that uses the available angles, with no range information involved, to drive the robot to the goal. Therefore, our method takes advantage of the strengths of omnidirectional vision, which provides a wide field of view and very precise angular information. We present a formal proof of the stability of the proposed control law. The performance of our approach is illustrated through simulations and different sets of experiments with real images.

## I. INTRODUCTION

### A. Problem Statement

Vision sensors have long been used to perform robot navigation [1] due to the high amount of information provided by cameras and their relatively low cost. Many research efforts in today's robotics continue to tackle the problem of visual control [2], [3]. A task of recognized interest in the vision-based navigation and control fields is *homing*, whose objective is to enable a robot to return to a previously visited location in an environment using visual information. Visual homing is often inspired by the mechanisms that certain animal species, such as insects, utilize to return to their known home location [4], [5]. A wide variety of methods have been employed to perform this task. For example, there are works on visual homing based on image distance [6], in scale space [7] or using frequency components [8].

Our approach is firstly motivated by the properties of omnidirectional vision sensors. They offer a wide field of view, which is a very interesting quality for navigation, and provide very precise angular information. In contrast, their radial information is strongly affected by distortion. In addition to this, it is known that the robustness of the performance of vision-based tasks can be improved if, instead of using the image information directly (either through appearance-based measurements or extracted features), the geometric models

that relate the views of a scene are employed. Thus, we propose a homing method for a robot moving on the ground plane that makes use of the angles between omnidirectional views extracted by means of the 1D trifocal tensor model. Our approach, which takes advantage of the planar motion constraint through the 1D tensor, employs only the angular visual information provided by omnidirectional images to obtain the relative angles between the view from the current position and a set of previously acquired reference views taken at different locations, any of which can be selected as the home (or goal) position. A stable control law based on the estimated angles is used to guide the robot to the target.

### B. Related Work

The interesting properties of omnidirectional cameras have motivated the proposal of a number of angle-based homing methods, being [9] an early work and [10], [11] more recent contributions. These approaches are purely feature-based: landmark angles in the images are used to generate a homing vector, which defines the direction of motion towards the goal. An alternative which is known to increase the robustness in the presence of wrong feature matches is to employ the models that express the intrinsic geometric relations between the views of a scene. A number of visual control methods have been presented using two-view models such as the epipolar geometry [12]–[15] and the homography [16]–[18].

The trifocal tensor is the model that encapsulates the geometric constraints of three views of a scene [19]. This model has been used for control purposes [20], [21]. In particular, robot navigation on the ground plane lends itself to the use of the 1D trifocal tensor, the matching constraint between three 1D views which minimally parameterizes them [22]. The 1D trifocal tensor, which can be computed linearly from point correspondences, allows to perform 2D projective reconstruction. In robotics, it has been used for 2D localization tasks [23], [24] and control [25], [26].

When compared with other approaches in the literature of visual navigation, ours is neither a pure topological map-based nor image memory-based approach, but it shares elements with both kinds of techniques. As the representation of the robot's environment necessary to perform navigation, we use a set of omnidirectional images. These serve us not only to build a connectivity graph of the environment, as in [27], [28], but also to store further information, which

This work was supported by Ministerio de Ciencia e Innovación/European Union (projects DPI2009-08126 and DPI2012-32100), by Ministerio de Educación under FPU grant AP2009-3430, and by DGA-FSE (group T04).

fits the *view graph* concept of [29]. We use the extracted geometric information relating the views to find connections between all the available reference images. Thus, in contrast with [27], [28], the graph we construct is not appearance-based, but as close to complete as possible. In addition to the graph, we also store the geometric information needed to perform homing between any two connected nodes in it using our proposed angle-based control law.

Some works in the literature use an image memory in a visual path following scheme [10], [17], [30]. We also employ a memory of stored images, which in our case serves as a representation of the environment that provides the references necessary to compute the geometric information that permits navigation. However, in contrast with these methods, our approach allows direct homing between two distant positions without having the need to travel along a rigid visual route of intermediate images. This long-range ability of our technique also differentiates it from local homing approaches, which are typically short-range [4]–[8], [12].

### C. Contributions relative to the literature

We believe that an important contribution of our work which sets it apart from the other visual homing approaches is the use of the 1D trifocal tensor to extract the geometric information. In addition to being particularly appropriate for the characteristics of omnidirectional cameras and planar motion, as discussed above, this tensor proves to be a strong constraint which allows the extraction of angular information in a remarkably robust way. The reasons for this are the presence of three views in the model and the possibility to perform triangle consistency checks to validate the estimations. This increases the robustness with respect to purely feature-based methods and approaches employing two-view geometric constraints. In addition, the angle-based control law we present has strong stability properties, further contributing to the robustness of the overall system.

We also present as contribution a method to compute the relative angles between all the available omnidirectional views. This is achieved through the use of the computation of the 1D trifocal tensor, a new angular disambiguation procedure and a novel approach for the indirect computation of angles. This method is what endows our technique with long-range direct homing capabilities.

The visual control method proposed in this paper was first presented in [25]. In the present paper, we provide several new contributions with respect to that previous work, namely:

- 1) Definition and analysis of the geometric visual three-view reconstruction ambiguities associated to our method.
- 2) Formal description of some aspects of the technique such as the connectivity of the reference views.
- 3) Stronger and more complete stability results. Namely, we provide in this paper both global asymptotic stability and local exponential stability proofs for our control law. In addition, the original control law has been slightly modified to make the motion of the robot smoother.

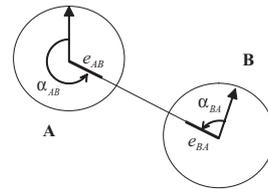


Fig. 1. Nomenclature and conventions used throughout the paper.  $e_{AB}$  is the epipole of view B in view A.  $\alpha_{AB}$  is the angle or direction of that epipole, i.e. its angular polar coordinate in view A. The reference axis for the measurement of the angles is given by the robot's direction of motion. The angles are measured counterclockwise between  $-\pi$  and  $\pi$ .

4) Extended discussion of the method and its characteristics, including a broader literature survey, in order to highlight the relevance of the presented work.

5) Two new sets of experiments, using images acquired both indoors and outdoors, which validate the performance of the approach and illustrate its robustness under varying working conditions.

The contents of the paper are organized as follows: Section II presents the procedure for the computation of all the angular information needed for the homing task. In section III the designed control strategy is described and its stability analysis is presented. Section IV provides an account of the results of the simulations and experiments conducted with real images. We discuss in section V a number of characteristics of the method and compare it with existing works. Finally, in section VI the conclusion and directions for future work are given. Some nomenclature and conventions used throughout the paper are illustrated in Fig. 1.

## II. COMPUTATION OF THE REFERENCE SET ANGLES

From the set of omnidirectional images we use to represent the environment, we can define an undirected graph  $G(V, E)$ . The nodes or vertices which make up the set  $V$  are the locations on the plane where each of the images were acquired. A link exists between two given nodes when the relative angles between their corresponding associated views are known, and  $E$  is the set of these links. Therefore, we can think of this graph as a topological map (i.e. a representation of an environment as a group of places and a group of connections between them). The adjacency matrix ( $A$ ) of  $G$  is defined such that  $A(i, j) = 1$  when a link in  $E$  exists between nodes  $i$  and  $j$  in  $V$ . Otherwise,  $A(i, j) = 0$ . In addition to the graph description, which expresses the connectivity of the nodes in our system, we must also store the relative angles between the views. We do so using a data structure  $M$  such that  $M(i, j) = \alpha_{ij}$  is the angle of the epipole of view  $j$  in view  $i$ .

Relevant features are extracted and matched between pairs of images on the reference set, and the resulting point correspondences are stored. We then start an estimation procedure that operates as follows:

- A set of four images (which we can call A, B, C and D) taken in two groups of three (e.g. A-B-C and B-C-D) are processed in each step. For each trio of images we obtain three-view point correspondences by taking the

common two-view matches between them. From a minimum number of seven point matches between the three views in each group, we can calculate two 1D trifocal tensors, and we can eventually obtain the angles of the epipoles in each of the views of the four-image set (as will be described in the following subsections). This way, the four associated graph nodes become adjacent.

- We run through the complete set of reference images calculating trifocal tensors and estimating the angles between the views. Whenever there is more than one estimation of a certain angle available, we simply choose the result that was obtained from the largest set of point matches. A possible alternative would be to compute a weighted average. After this stage is completed, we usually end up with an incomplete graph.

- We fill the graph's adjacency matrix using the indirect estimation procedure that will be described in section II-B. This procedure is aimed at computing the unknown angles between the views from the already known values. For every pair of nodes  $i, j$  such that  $A(i, j) = 0$ , we search for two other nodes  $k, l$  such that  $A(i, k), A(i, l), A(j, k), A(j, l)$  and  $A(k, l)$  all equal 1. When these conditions hold, we are able to compute the angles between  $i$  and  $j$ , and consequently make  $A(i, j) = 1$  and calculate the value  $M(i, j)$ .

The typical way in which the graph's adjacency matrix becomes filled is the following: when several locations are adjacent physically, the relative angles between them can be computed directly thanks to the large amount of common visual information between the views, and the graph nodes corresponding to these locations are also adjacent. Thus, initially we have the whole set of images divided in groups of connected locations (with adjacent associated graph nodes). If two of these groups are adjacent physically, it is very likely that they have at least two nodes in common. When this is the case, we can connect all the elements in the two groups using the indirect angle computation procedure, and all the nodes in the two groups become adjacent. Following this procedure, the adjacency matrix gets gradually filled until every pair of nodes are adjacent (i.e. the graph is complete) or it is not possible to find new connections between the nodes. Next, we describe the different steps of the process through which the relative angles between the views are obtained.

#### A. Angles from the 1D trifocal tensor

The trifocal tensor is the mathematical entity that encapsulates all the geometric relations between three views that are independent of scene structure. In particular, the 1D trifocal tensor relates three 1D views on a plane, and presents interesting properties; namely, it can be estimated linearly from a minimum of seven three-view point matches (or five, if the calibration of the cameras is known [31]), and 2D projective reconstruction can be obtained from it.

1) *Foundations: 1D trifocal tensor computation and epipole extraction:* The projections of a given point in three 1D views (which we will refer to as  $A, B$  and  $C$ ) on a plane are related by the following trilinear constraint [22]:

$$\sum_{i=1}^2 \sum_{j=1}^2 \sum_{k=1}^2 T_{ijk} u_i^A u_j^B u_k^C = 0, \quad (1)$$

where  $T_{ijk}$  are the elements of the *1D trifocal tensor*,  $\mathbf{T} \in \mathbb{R}^{2 \times 2 \times 2}$ , between the views, and  $\mathbf{u}^A, \mathbf{u}^B$  and  $\mathbf{u}^C$  are the homogeneous coordinates of the projections of the point in each view.  $\mathbf{T}$  is defined up to a scale factor and therefore can be calculated, in the uncalibrated case, from a minimum set of seven point correspondences across the views.

The process we follow to estimate  $\mathbf{T}$  starts by detecting relevant features in three omnidirectional images, e.g. by means of the SIFT keypoint extractor [32], and finding matches between them. The angles ( $\alpha$ ) of the matched image points, measured counterclockwise from the vertical axis, are converted to a 1D projective formulation, with the corresponding homogeneous 1D coordinates being  $(\sin \alpha, \cos \alpha)^T$ . In this mapping, the distinction between angles differing by  $\pi$  is lost.

Each of the point matches in 1D projective coordinates gives rise to an equation of the form of (1). If more than seven correspondences are available, we find a least squares solution to the resulting system of linear equations through Singular Value Decomposition [19], [24]. In this process, a robust estimation method (RANSAC) is employed in order to reject wrong matches.

After  $\mathbf{T}$  has been estimated, we extract the epipoles from it using a procedure taken from [23], [33] that we briefly describe next. A 1D homography is a mapping between projected points in two lines (two of the 1D views, in our case) induced by another line. From the coefficients of the trifocal tensor, we can directly extract what are known as the *intrinsic homographies*. For example, the two intrinsic homographies from  $A$  to  $B$ ,  $\mathbf{K}_{AB}$  and  $\mathbf{L}_{AB}$ , are obtained by substituting the lines defined by  $\mathbf{u}^C = (1, 0)^T$  and  $\mathbf{u}^C = (0, 1)^T$  in (1), yielding

$$\mathbf{K}_{AB} = \begin{bmatrix} -T_{211} & -T_{221} \\ T_{111} & T_{121} \end{bmatrix}, \mathbf{L}_{AB} = \begin{bmatrix} -T_{212} & -T_{222} \\ T_{112} & T_{122} \end{bmatrix}. \quad (2)$$

Now,  $\mathbf{H}_A = \mathbf{K}_{AB} \mathbf{L}_{AB}^{-1}$  is a homography from  $A$  to itself; by definition, the epipoles are the only points that are mapped to themselves by such a homography, i.e.:  $\mathbf{e}_{AB} = \mathbf{H}_A \mathbf{e}_{AB}$  and  $\mathbf{e}_{AC} = \mathbf{H}_A \mathbf{e}_{AC}$ . Therefore we can calculate them as the eigenvectors of matrix  $\mathbf{H}_A$ . It is important to note, though, that with this method we do not know which of the other two views ( $B$  or  $C$ ) each of the two recovered epipoles corresponds to. By mapping this pair of epipoles to the other views through the intrinsic homographies, we finally obtain the six epipoles of the set of three 1D views.

2) *Ambiguity resolution:* There are three ambiguities that need to be resolved in order to determine the correct values of the angles of the 2D epipoles from the values of the epipoles extracted using the 1D trifocal tensor.

Firstly, as mentioned in section II-A.1, an epipole in a given view recovered from the 1D trifocal tensor may be assigned to any of the two other views. This results in

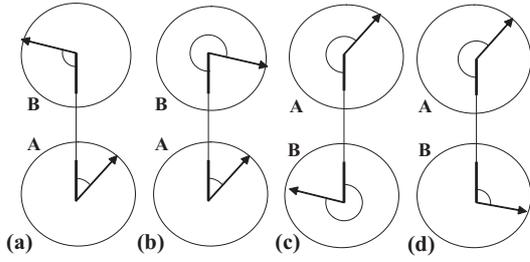


Fig. 2. Four possible 2D reconstructions from the epipoles between two views extracted from the 1D trifocal tensor (top). The relations between the angles of the projections of matched points (e.g.  $P_1$  and  $P_2$ ) in two aligned views can be used to resolve the 2D reconstruction ambiguities (below).

two possible solutions in the assignment of the set of six epipoles between the three views. As shown in [22], [31], both solutions give completely self-consistent 2D projective reconstructions, regardless of the number of point matches between the views. This fundamental ambiguity in the 2D reconstruction from three 1D views can only be resolved through the use of a fourth view, as noted in [23]. We propose a new method to resolve the ambiguity that operates in the following way: having a group of four views (which we can call A, B, C and D), we calculate two different trifocal tensors between them; for instance, the tensor relating A, B and C, and the tensor between B, C, and D. Since the epipoles between B and C must be identical in the two estimations, by detecting the common (or, in a real situation, the closest) epipoles in these two views we can disambiguate the assignment of the complete set of epipoles.

The origin of the two other ambiguities lies in the fact that in the mapping of 2D image points to 1D projective coordinates, the distinction between bearings differing by  $\pi$  is lost. The angle of a recovered 1D epipole  $(e_1, e_2)^T$  is obtained as  $\arctan(e_1/e_2)$  in 2D. As a consequence, from the 1D epipole we can extract two different angles in a 2D view, separated by  $\pi$  radians. There are, therefore, four possible solutions for the values of the epipoles between two given views A and B, which may be interpreted as emerging from two combined reconstruction ambiguities; namely, an ambiguity in the direction of the translation vector from view A to view B, which accounts for the difference between solutions (a) and (c) in Fig. 2, and an ambiguity of  $\pi$  radians in the orientation of view B, illustrated, for example, by solutions (a) and (b) in the same figure.

This double ambiguity for a set of two views might be resolved through point reconstruction, but instead we propose a simple method employing only the angles of matched image points. The method takes advantage of the properties

TABLE I  
DISAMBIGUATION OF THE ANGLES OF THE EPIPOLES IN TWO VIEWS

$test_1$	$test_2$	$\alpha_{AB}$	$\alpha_{BA}$	Case in Fig. 2
1	1	$\alpha_{AB}^s$	$\alpha_{BA}^s$	(a)
0	1	$\alpha_{AB}^s$	$\alpha_{BA}^s + \pi$	(b)
1	0	$\alpha_{AB}^s + \pi$	$\alpha_{BA}^s + \pi$	(c)
0	0	$\alpha_{AB}^s + \pi$	$\alpha_{BA}^s$	(d)

of the optical flow between two omnidirectional images when they are aligned. To do so, it exploits the relative angular differences between the projections of a scene point in the two images. This is illustrated in Fig. 2. The complete procedure for the computation of the relative angles between views from the 1D trifocal tensor, resolving all the existing ambiguities, is provided in algorithm 1.

Next, we determine what the possible reconstructions are for a trio of 1D views in the plane, starting from the knowledge of the 1D epipoles between them. As we have already seen, a 1D epipole can represent two different angles in the plane separated by  $\pi$  radians, which gives four possible reconstructions of the 2D relative geometry between two views. With three views, we have six epipoles, and the combination of all the possible values of the angles gives  $2^6 = 64$  possible reconstructions of the geometry of the views. However, by considering geometric consistency constraints between the three views, it can be shown that only  $2^3 = 8$  of these possibilities are actually valid, jointly coherent reconstructions.

When we convert the points in three omnidirectional images with a common image plane into 1D coordinates, the situation is equivalent to having three 1D cameras in a common 2D space. We can express the relative locations of these cameras with a rotation matrix  $\mathbf{R} \in \mathbb{R}^{2 \times 2}$  and a translation vector  $\mathbf{t} \in \mathbb{R}^{2 \times 1}$ . By considering one of the cameras ( $\mathbf{P}_A$ ) as the fixed origin of the reference system, we have the following projection matrices:  $\mathbf{P}_A = [\mathbf{I} | \mathbf{0}]$ ,  $\mathbf{P}_B = [\mathbf{R}_B | \mathbf{t}_B]$  and  $\mathbf{P}_C = [\mathbf{R}_C | \mathbf{t}_C]$ . The rotation matrices are as follows:

$$\mathbf{R}_B = \begin{bmatrix} \cos \phi_B^s & \sin \phi_B^s \\ -\sin \phi_B^s & \cos \phi_B^s \end{bmatrix}, \quad \mathbf{R}_C = \begin{bmatrix} \cos \phi_C^s & \sin \phi_C^s \\ -\sin \phi_C^s & \cos \phi_C^s \end{bmatrix}, \quad (3)$$

where we can express:  $\phi_B^s = \pi - \alpha_{AB}^s + \alpha_{BA}^s$  and  $\phi_C^s = \pi - \alpha_{AC}^s + \alpha_{CA}^s$ . The superscript  $s$  alludes to angles extracted directly from the 1D epipoles arbitrarily (i.e. without having been disambiguated). The translation vectors can also be worked out:

$$\mathbf{t}_B = \begin{bmatrix} -e_{AB} \cos \phi_B^s - \sin \phi_B^s \\ e_{AB} \sin \phi_B^s - \cos \phi_B^s \end{bmatrix} \quad (4)$$

$$\mathbf{t}_C = \begin{bmatrix} -e_{AC} \cos \phi_C^s - \sin \phi_C^s \\ e_{AC} \cos \phi_C^s - \cos \phi_C^s \end{bmatrix}, \quad (5)$$

where  $e_{AB}$  and  $e_{AC}$  are the epipoles in inhomogeneous coordinates. The vectors  $\mathbf{t}_B$  and  $\mathbf{t}_C$  are defined up to scale. Actually, there is only one scale in the three-view

---

**Algorithm 1** Disambiguation of the relative angles between views A,B,C and B,C,D from the 1D epipoles
 

---

- 1) Let us consider the 1D epipoles in inhomogeneous coordinates in view B:  $e_{B1}, e_{B2}$  extracted from the tensor computed between A, B and C, and  $e'_{B1}, e'_{B2}$  extracted from the tensor computed between B, C and D. Find  $i \in \{1, 2\}, j \in \{1, 2\}$  such that  $\min\{|\arctan(e_{Bi}) - \arctan(e'_{Bj})|, |\arctan(e_{Bi}) - \arctan(e'_{Bj}) - \pi|\}$  is minimum. Then compute  $\alpha_{BC} = (\arctan(e_{Bi}) + \arctan(e'_{Bj}))/2$ , and assign the angles of all the other epipoles accordingly. The angles have an ambiguity of  $\pi$ .
  - 2) For each pair of images (let us take A and B as example):
    - a) Define  $\alpha_{AB}^s = \{\alpha_{AB} \vee \alpha_{AB} + \pi\}$  and  $\alpha_{BA}^s = \{\alpha_{BA} \vee \alpha_{BA} + \pi\}$  so that both  $\alpha_{AB}^s$  and  $\alpha_{BA}^s \in (0, \pi]$
    - b) Rotate all the  $m$  points matched between A and B ( $\alpha_{Ak}$  and  $\alpha_{Bk}$ , measured counterclockwise):  
For  $k = 1, \dots, m$ ,  $\alpha_{Akr} = \alpha_{Ak} - \alpha_{AB}^s$ ,  $\alpha_{Bkr} = \alpha_{Bk} - \alpha_{BA}^s$
    - c) Test 1: Relative orientation of the two aligned images. For  $k = 1, \dots, m$ ,  
 $nr(k) = \text{sign}(\alpha_{Akr} \cdot \alpha_{Bkr}) \cdot \min(|\sin \alpha_{Akr}|, |\sin \alpha_{Bkr}|)$ . If  $\sum_{k=1}^m nr(k) > 0$ , then  $test_1 = 0$ , otherwise  $test_1 = 1$
    - d) Test 2: Sign of translation from A to B. If  $test_1 = 1$ , then for  $k = 1, \dots, m$ ,  $\alpha_{Bkr} = \alpha_{Bkr} + \pi$ .  
For  $k = 1, \dots, m$ , if  $\text{sign}(\alpha_{Akr} \cdot \alpha_{Bkr}) = 1$  then  $st(k) = \text{sign}(|\alpha_{Akr}| - |\alpha_{Bkr}|) \cdot (|\alpha_{Akr}| - |\alpha_{Bkr}|)^2$ ,  
otherwise  $st(k) = 0$ . If  $\sum_{k=1}^m st(k) > 0$ , then  $test_2 = 0$ , otherwise  $test_2 = 1$
    - e) Obtain from Table I the unambiguous angles  $\alpha_{AB}$  and  $\alpha_{BA}$
  - 3) Check the joint coherence of the disambiguations of the pairs of views taken in groups of three, by verifying in Table II that the obtained three-view reconstructions are valid.
- 

TABLE II

THE EIGHT POSSIBLE RECONSTRUCTIONS OF THE CAMERA LOCATIONS FROM THE EPIPOLES IN THREE VIEWS

$P_B =$	$\begin{bmatrix} R_B \\ t_B \end{bmatrix}$	$\begin{bmatrix} R_B \\ t_B \end{bmatrix}$	$\begin{bmatrix} -R_B \\ -t_B \end{bmatrix}$	$\begin{bmatrix} -R_B \\ -t_B \end{bmatrix}$
$P_C =$	$\begin{bmatrix} R_C \\ t_C \end{bmatrix}$	$\begin{bmatrix} -R_C \\ -t_C \end{bmatrix}$	$\begin{bmatrix} R_C \\ t_C \end{bmatrix}$	$\begin{bmatrix} -R_C \\ -t_C \end{bmatrix}$
$P_B =$	$\begin{bmatrix} R_B \\ -t_B \end{bmatrix}$	$\begin{bmatrix} R_B \\ -t_B \end{bmatrix}$	$\begin{bmatrix} -R_B \\ t_B \end{bmatrix}$	$\begin{bmatrix} -R_B \\ t_B \end{bmatrix}$
$P_C =$	$\begin{bmatrix} R_C \\ -t_C \end{bmatrix}$	$\begin{bmatrix} -R_C \\ t_C \end{bmatrix}$	$\begin{bmatrix} R_C \\ -t_C \end{bmatrix}$	$\begin{bmatrix} -R_C \\ t_C \end{bmatrix}$

reconstruction, since a relation exists between the magnitudes of  $t_B$  and  $t_C$  through the triangle formed by the views:

$$\frac{\|t_C\|}{\|t_B\|} = \left| \frac{\sin(\alpha_{BA}^s - \alpha_{BC}^s)}{\sin(\alpha_{CA}^s - \alpha_{CB}^s)} \right|. \quad (6)$$

The eight possible reconstructions, up to scale, from the epipoles between three views are shown in Table II.

Thus, after using the proposed two-view disambiguation method, we check if the views taken in groups of three give jointly coherent reconstructions. In addition, we also check that the triangle formed by the locations of the three views is consistent (i.e. that the sum of its angles is sufficiently close to  $\pi$ ). By doing so the angles between every trio of views are estimated robustly.

### B. Complete solution of four-view sets

In practice, it is usually not possible to find matches across all the images. Next, we propose a method to compute all the angular information using the matches between sets of adjacent or close images. A geometric setting of the type shown in Fig. 3, where two triangles are known between the locations of four views, comes up in our method every time we estimate two trifocal tensors from a four-view set. This section describes the method employed to calculate the two unknown angles in this configuration.

We use the notation  $\widehat{ABC}$  to allude to the values ( $> 0$ ) of the angles in a triangle. Without loss of generality, we

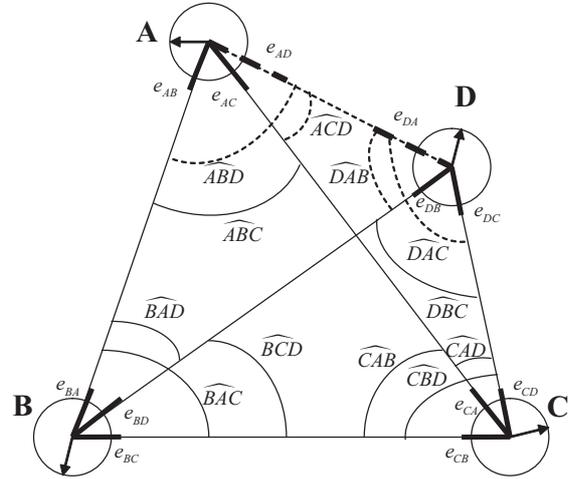


Fig. 3. Geometric setting with four views and two known triangles.

can formulate the problem in the following terms: all the angles from every view to the others in the set are known except the angles of the epipoles between views A and D. Therefore all the angles in the four triangles formed by the set of four views are known, except the ones including both vertices A and D (represented with dashed lines in Fig. 3). Our objective is to calculate the angles  $\alpha_{AD}$  and  $\alpha_{DA}$  of the epipoles  $e_{AD}$  and  $e_{DA}$ , which can be directly obtained from the knowledge of the angles of the triangles at those vertices. We start by applying the law of sines on the set of four triangles (ABC, ABD, ACD and BCD), which finally yields the following expression

$$\frac{\sin \widehat{ABD}}{\sin \widehat{ACD}} = K_A, \quad (7)$$

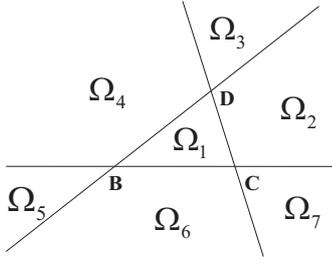


Fig. 4. Seven regions where point A can be located.

TABLE III

VALUES OF SIGNS FOR THE DIFFERENT REGIONS IN WHICH A MAY LIE

Region of vertex A	Relation between angles at vertex A	$sign_1$	$sign_2$
$\Omega_1$	$\widehat{ACD} = 2\pi - \widehat{ABD} - \widehat{ABC}$	-1	1
$\Omega_2, \Omega_5$	$\widehat{ACD} = \widehat{ABD} + \widehat{ABC}$	1	-1
$\Omega_3, \Omega_6$	$\widehat{ACD} = \widehat{ABC} - \widehat{ABD}$	1	1
$\Omega_4, \Omega_7$	$\widehat{ACD} = \widehat{ABD} - \widehat{ABC}$	-1	-1

where  $K_A$  is a known value given by

$$K_A = \frac{\sin \widehat{CBD} \cdot \sin \widehat{BAD}}{\sin \widehat{BCD} \cdot \sin \widehat{CAD}}. \quad (8)$$

Using the intrinsic relationship between the three angles at vertex A and applying trigonometric identities, we can calculate the individual values of the angles in (7). We must, however, take into account the fact that the location of A with respect to the other three vertices changes the geometry of the set and, consequently, the relation between the angles at the aforementioned vertex. Therefore, we need to divide the plane into seven regions, as shown in Fig. 4, to account for these differences. It turns out that the expression that gives the angle  $\widehat{ABD}$  has the same form in all cases (i.e. for all regions), but the signs of two of its terms, denoted as  $sign_1$  and  $sign_2$ , are dependent on the region where A lies

$$\widehat{ABD} = \arctan \frac{sign_1 \cdot K_A \sin(\widehat{ABC})}{1 + sign_2 \cdot K_A \cos(\widehat{ABC})}. \quad (9)$$

We can easily determine the region in which A is located using the known angles of the epipoles in views B and C, and choose the appropriate values of  $sign_1$  and  $sign_2$  as shown in table III.

The angle of the epipole of view D in view A is finally obtained as follows

$$\alpha_{AD} = \begin{cases} \alpha_{AB} + \widehat{ABD}, & \text{if } 0 \leq \alpha_{BA} - \alpha_{BD} < \pi \\ \alpha_{AB} - \widehat{ABD}, & \text{if } 0 > \alpha_{BA} - \alpha_{BD} \geq -\pi \end{cases}. \quad (10)$$

The angle of the epipole in view D of view A ( $\alpha_{DA}$ ) can be calculated through a completely analogous process, simply interchanging the roles of vertices A and D. The results are validated using geometric consistency checks. By employing the procedure we have just presented, we can calculate the two unknown angles and thus obtain the complete set of angles between the four views. In addition, this method is useful for two other purposes within our homing technique:

- In the initial stage, described in the beginning of section II, this method allows to fill in the missing elements in the matrix of epipole angles, corresponding to pairs of views that could not be linked directly due to the impossibility to find a sufficiently large set of three-view matches between them.

- During homing, it enables us to obtain all the angles needed to generate the motion commands employing a minimum number of three views; we only need to compute the trifocal tensor between the current image taken by the robot and two of the reference images, which reduces the cost of the algorithm.

### III. HOMING STRATEGY

In this section we describe the strategy designed in order for the mobile robot to perform homing. We assume the robot moves on the ground plane and has unicycle kinematics. The homing method is based solely on the computation of the angles between a series of omnidirectional views of the environment. This group of snapshots consists of the image taken by the robot from its current position and a set of previously acquired reference images, which includes an image obtained at the desired target location. The angles between the views on the reference set have been previously computed and stored, as described in section II. Therefore, only the angles between the current and the reference views must be worked out during homing.

In every step of the robot's motion, the camera takes an omnidirectional image, from which key points are extracted. When sufficient point matches are found between the current and two of the reference images, the 1D trifocal tensor is calculated as detailed in section II-A.1. From the tensor, aided by the knowledge of the angles on the reference set, we can extract the angles between the current and the two other views. Finally, with the method explained in section II-B all the angles of the epipoles in all the views can be computed.

#### A. Control law

For every reference view  $R_i(x_i, z_i, \varphi_i)$  (where  $x_i, z_i$  and  $\varphi_i$  define its position and orientation in the ground plane), the difference between the angles of its epipoles with respect to the current and goal views defines an angular sector of size  $S_i = |\alpha_{iC} - \alpha_{iG}|$ , as illustrated in Fig. 5. We use the average value of the angular sizes of these sectors to set the linear velocity at which the robot will move toward the target position

$$v = k_v \cos \alpha_{CG} \cdot \frac{1}{n} \sum_{i=1}^n S_i, \quad (11)$$

where  $k_v > 0$  is a control gain. As the robot moves closer to the goal, the mean size of the angular sectors seen from the reference positions will become smaller; thus, the robot's linear velocity will gradually decrease and eventually become zero when the target is reached. The cosine term in (11) ensures that  $v$  has the right sign; when the target is behind the robot,  $\cos \alpha_{CG}$  will be negative, therefore generating backward motion. In addition, the cosine improves the behavior by gradually reducing the vehicle's translational speed when it is not pointing to the goal.

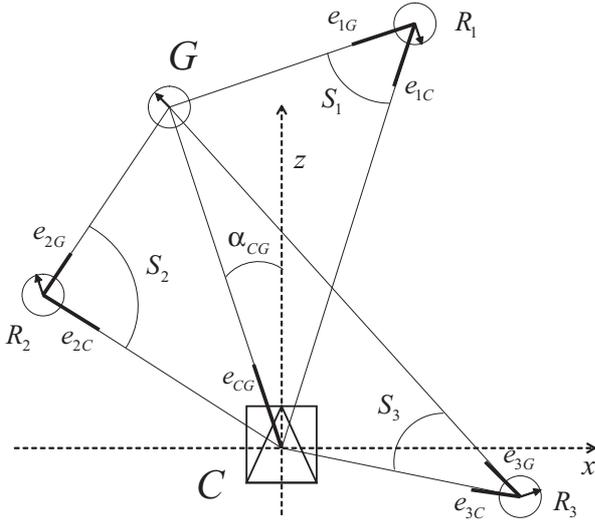


Fig. 5. Elements involved and angles employed in the homing strategy.  $C$  is the robot's current location, at the coordinate origin.  $G$  is the goal location.  $R_i$  are reference views. Three of the  $n$  views on the reference set are depicted as example.

The direction in which the robot travels is determined by the angle at which the goal position is seen from the current location, i.e. the angle  $\alpha_{CG}$  of the epipole  $e_{CG}$ . The angular velocity of the control law is given by

$$\omega = k_\omega(\alpha_{CG} - \alpha_{CG}^d), \quad (12)$$

$$\alpha_{CG}^d = \begin{cases} 0 & \text{if } |\alpha_{CG_0}| \leq \frac{\pi}{2} \\ \pi & \text{if } |\alpha_{CG_0}| > \frac{\pi}{2} \end{cases},$$

where  $k_\omega > 0$  is a control gain, and  $\alpha_{CG_0}$  is the value of  $\alpha_{CG}$  at the start of the execution. From a minimum number of four reference views, one of which would be the view from the target location, the robot will navigate to the home position. Note that the orientation in which the robot reaches the target position is not controlled, since, by definition, the purpose of the homing task is getting to the goal location. In addition, final orientation correction is less relevant when one uses robots equipped with omnidirectional cameras, since they always provide a full view of the environment regardless of the robot's orientation.

### B. Stability Analysis

We first define  $d_i$  as the distance between reference view  $i$  and the goal position (i.e. the length of line segments  $\overline{R_i G}$  in Fig. 5), and  $d_{min}$  as the minimum of such distances.

As shown in theorem 1 (section A of the appendix), the system under the proposed control law (11), (12) is globally asymptotically stable if  $k_\omega > k_v \cdot \pi / d_{min}$ . This result means that we can ensure the stability of the system if we have an estimate of the value of  $d_{min}$  for the particular set of reference views we are working with. This estimate does not need to be precise, since what we have found is a fairly conservative bound for the value of  $k_\omega$ . In practice, the

system will be stable with angular velocity gains lower than this threshold.

Additionally, as shown in proposition 1 (section B of the appendix), the system under the proposed control law (11), (12) is locally exponentially stable. The actual value of the exponential decay parameter (which we call  $\lambda_{min}$ ) in a particular case will depend on both the geometric distribution of the reference views and the trajectory along which the robot approaches the goal. The less collinear the reference views are, the faster the system will be guaranteed to converge.

## IV. EXPERIMENTAL RESULTS

The performance of the proposed method has been tested both in simulation and with real images.

### A. Simulations

For the first simulation we present, the reference views were positioned forming a square grid. A randomly distributed cloud of 200 points in 3D was generated and projected in each camera. Four sample homing trajectories with a 25-view reference set and the evolutions of their corresponding motion commands are displayed in Fig. 6. The cloud of points is also shown. One of the trajectories starts from a position outside of the grid of reference views. As can be seen, it behaves in the same way as the other three. In all four trajectories the motion is smooth and the robot converges to the goal position.

The reference views can be laid out in any arbitrary configuration (as long as sufficient geometric diversity on the plane is guaranteed). We illustrate this fact with the simulation shown in Fig. 7. Eight reference images are used, and they are placed at arbitrary locations in the environment. The plot shows the paths from a common initial location (marked with a square) to the positions associated to the reference images, i.e. a different image is selected as the target each time.

We also added Gaussian noise to the angles of the projected points to evaluate the performance of the homing method. Fig. 8 displays the final position error obtained after adding variable noise in simulations with sets of 4 (the minimum number for our method), 8 and 16 reference images. Increasing the number of reference views makes the system more robust to noise, since the linear velocity of the control is computed by averaging out the contributions of all the views.

### B. Experiments with real images

In order to assess the performance of our homing method, we tested it with three diverse sets of real omnidirectional images.

1) *First indoor experiment:* The images for this first experiment were obtained in a laboratory setting. The experimental setup is illustrated in Fig. 9. It consisted of an ActivMedia Pioneer nonholonomic unicycle robot base with a catadioptric vision system, made up of a Point Grey FL2-08S2C camera and a Neovision HS3 hyperbolic mirror,

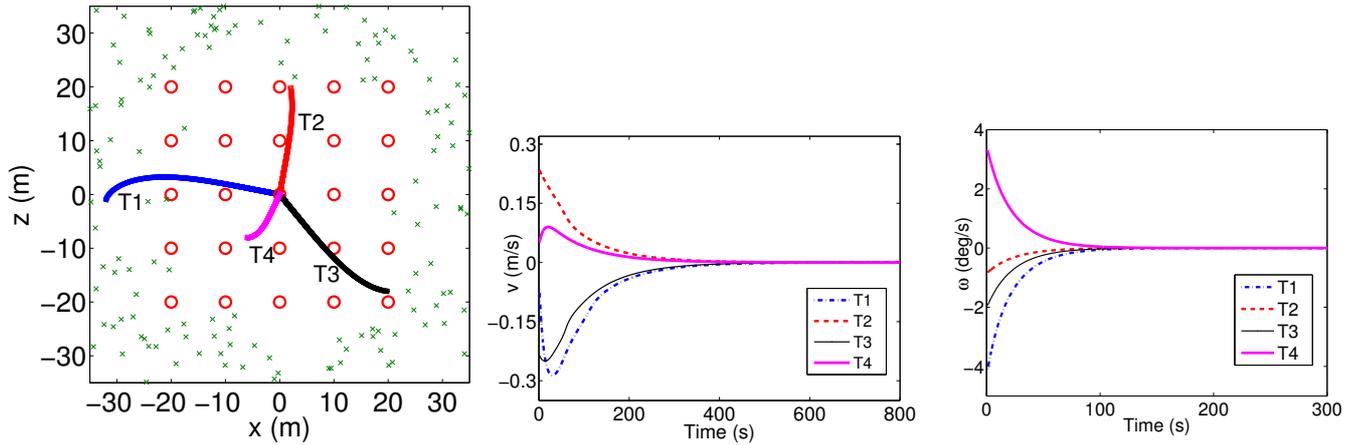


Fig. 6. Simulation results. Left: robot paths for four different trajectories with a reference set consisting of 25 images acquired in positions forming a grid. The goal view is at position (0,0) for all the trajectories. The locations of the reference images are marked as circles. The cloud of 200 points used to compute the 1D trifocal tensors is also shown. Center: Linear velocity for the four trajectories. Right: angular velocity for the four trajectories.

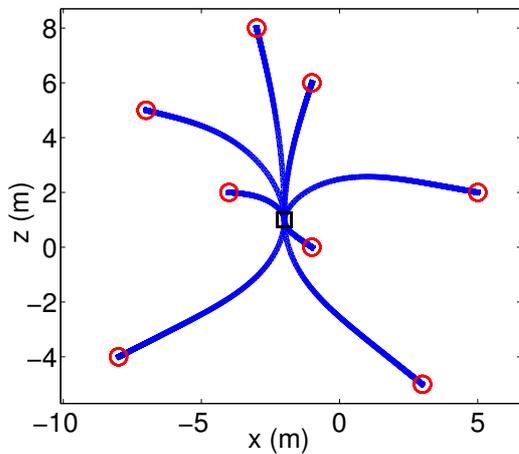


Fig. 7. Simulation results: robot paths with 8 reference images in arbitrary locations. The different paths obtained by taking each of the reference images (marked as circles) as the goal, from a common starting location (marked as a square), are shown.

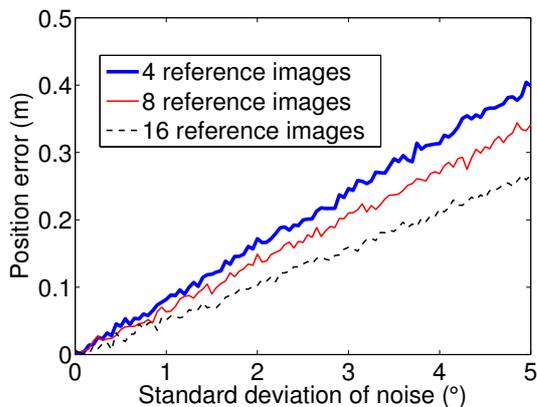


Fig. 8. Simulation results: final position error as a function of added Gaussian noise for reference sets of 4, 8 and 16 images.



Fig. 9. Omnidirectional camera (left) and complete setup (right) used for the first indoor experiment and the outdoor experiment.

mounted on top. The resolution of the employed images was  $800 \times 600$  pixels. The imaging system was used without specific calibration other than the assumption that the camera and mirror axis are vertically aligned.

To generate the reference set of views, 20 images were acquired from locations forming a  $5 \times 4$  rectangular grid with a spacing of 1.2 m., thus covering a total area of  $4.8 \times 3.6 \text{ m}^2$ . Features in the images were extracted and matched using SIFT, and a RANSAC robust estimation was used to calculate the 1D trifocal tensors between the views. The number of three-view correspondences employed to obtain the trifocal tensor estimations lied in the range of 30 to 80. Although images taken on opposite sides of the room could not be matched, the connections between adjacent or close sets of views were sufficient to recover the relative angles of the complete reference set. Vector field representations for two different goal locations within the grid are displayed in

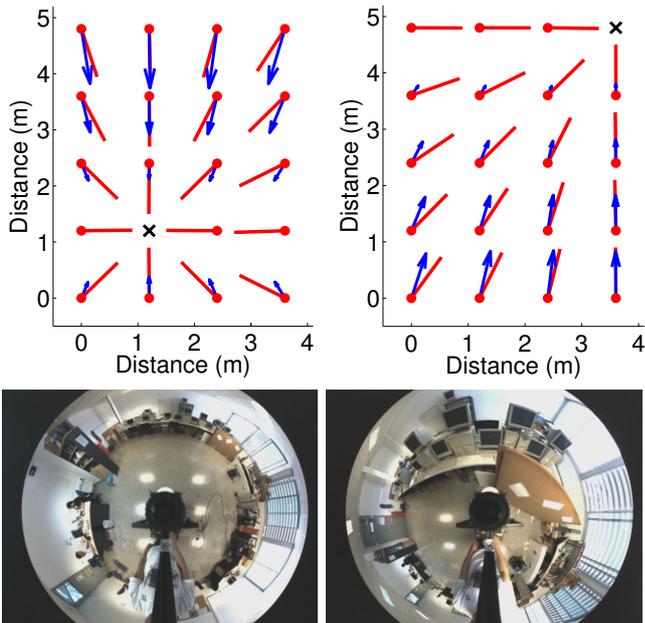


Fig. 10. Top: displacement vectors (arrows) and directions of the epipoles (line segments) with respect to the goal estimated at every reference position for two different goal locations (marked with a cross) in the real setting of the first indoor experiment. Bottom: goal images corresponding to the homing vector representations on top.

Fig. 10.

The arrows in this figure represent displacement vectors and have the following interpretation: if we consider a robot with unicycle kinematic constraints initially situated on each of the reference spots, with its initial orientation aligned with the  $Y$  axis of the plot, the arrows represent the displacement that the robot would perform according to the proposed control law. All the vectors have been scaled by an equal factor (the scale is different in the two plots, for better visualization). As can be seen, the magnitude of the vectors becomes larger as the distance to the target increases. Notice in both plots that in the positions from which the angle to the goal is  $\pm\pi/2$ , the robot does not translate initially. In these cases it first executes a rotation, due to the cosine term introduced in the linear velocity of the control law (11), before moving towards the goal. The line segments in Fig. 10 show the estimated directions of the epipoles of the goal position (i.e. the homing vectors) in each of the reference locations. The accuracy of the results obtained in this experiment is remarkable. The outliers in the set of putative point correspondences are rejected through the robust computation of the 1D trifocal tensor, and the availability of a fairly large number of good three-view matches makes it possible to compute the angles between views very precisely.

2) *Second indoor experiment:* In order to assess the performance of our method in more involved scenarios, we tested it on the dataset of omnidirectional images used in [34]. The images in this set (see examples in Fig. 11) were provided unwrapped and had a resolution of  $720 \times 120$  pixels. We note here that our approach can be applied

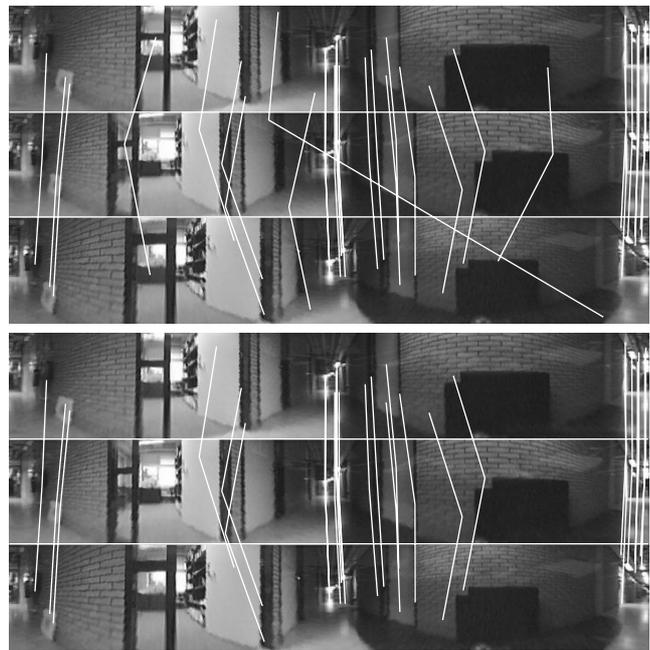


Fig. 11. Example of a trio of images from the data set of the second indoor experiment [34] with their putative SIFT correspondences joined by lines (top). Feature matches remaining after the robust computation of the 1D trifocal tensor (bottom).

indistinctly to both wrapped and unwrapped omnidirectional images. This data set covers an office environment comprising a long corridor (around 14 m. of length) and several rooms, and its images were acquired with separations of 0.3 or 0.5 m. The number of images used in this experiment was 213. Despite this denser coverage of the environment with respect to the first experiment, the extraction and matching of SIFT points between the views was much more complicated due to the lower quality of the images. Thus, it was necessary to select the parameters of the SIFT extractor appropriately, which was done following the guidelines described in [7]. A number of 12 three-view correspondences was found to be the minimum necessary for our method to operate on this data set. An example of the matching process for a trio of images from this set, showing the initial SIFT correspondences between them and the matches remaining after the robust computation of the 1D trifocal tensor, is displayed in Fig. 11.

The results of this experiment are illustrated in Fig. 12. The top three plots in the figure display the homing vectors from all the reference positions to the target location, for three goals chosen in different parts of the environment (inside a room, in the center of the corridor, and near the opposite end of it). As can be seen, the results are fairly accurate and robust despite the low quality of the images. This is also true when the positions are very far from the goal or even located in a different room. For a small group of locations in the set it was not possible to compute the angular information and resulting homing vector. These cases occurred when either the number of matched features was insufficient or incoherent results (which were automatically

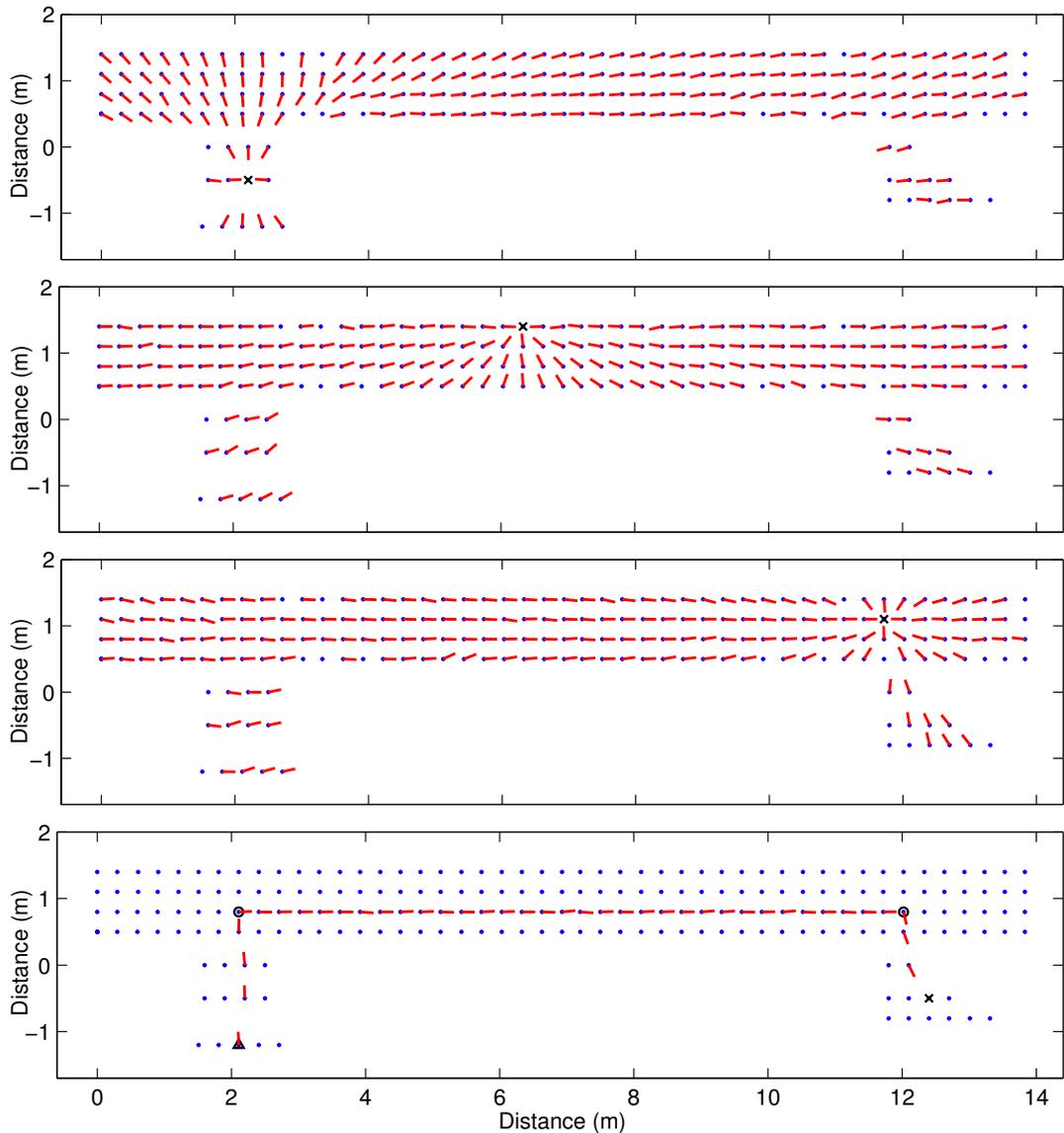


Fig. 12. Experimental results for the data set of the second indoor experiment. The environment consists of two rooms and a corridor. Top three plots: the computed homing vectors from all the positions of the reference views in the set to the goal location, marked with a cross, are shown for three different examples. Bottom plot: sampled path along a chain of reference positions from a starting location (triangle) to a homing destination (cross) passing through two homing sub-goals (circles).

detected and rejected by our method through geometric consistency checks) were obtained.

Obviously, it would be necessary to define intermediate goals in order for the robot to navigate in a setting of this kind. If, for instance, we would like to go from the room on the left (room 1) to the room on the right (room 2) in this environment, the sequence of steps to follow would be: first move through the door out of room 1 to a reachable goal position in the corridor, then travel to the other side of the corridor to a goal position in front of the door of room 2, then get into room 2 and reach the final goal. That is, a global homing task has to be divided in several *direct* homing steps, in each of which the goal location must be directly reachable from the starting position. We illustrate the results of an example of such a task in the bottom plot of Fig. 12.

The homing vectors along a sampled approximation of the path that the robot would follow are shown. The intermediate goals have been selected manually. Several aspects related with the long-range operation of our homing method are treated in the discussion section (V) of the paper.

We also present results with regard to the connectivity of the views in this experiment in Fig. 13. Connectivity graphs and their associated adjacency matrices are displayed. As can be seen, only images that are close physically can be initially connected. These links are obtained from the computation of the 1D trifocal tensor between sets of three views. Then, the matrix becomes filled using the indirect angle computation procedure. Two intermediate cases of this process are shown in Fig. 13, illustrating how the connections between views are progressively computed. Eventually, almost the complete

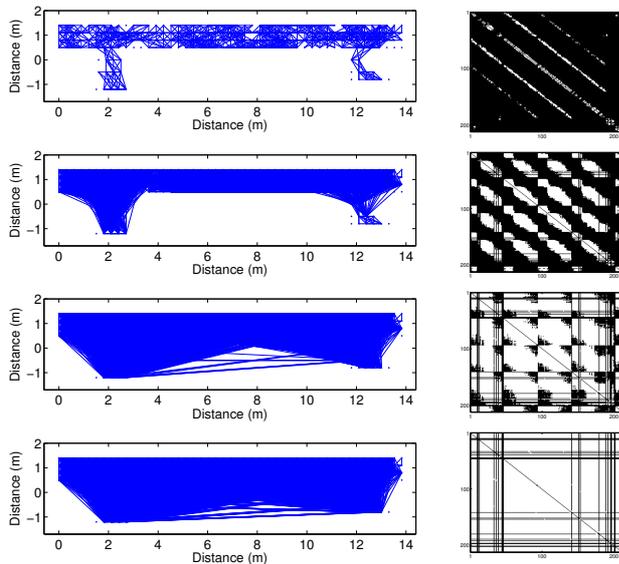


Fig. 13. Connectivity between views for the second indoor experiment. Left column: graph representations of the connectivity, with lines joining the connected locations in the environment after using our approach. Right column: color-coded adjacency matrices corresponding to the connectivity graphs on the left. White color means '1', black color means '0'. The values in the axes denote the indices of the images, ranging from 1 to the number of available images. From top to bottom, four cases are shown: the initial connectivity, two intermediate cases, and the final connectivity results.

set becomes interconnected, as shown by the plots in the bottom row of the figure. The rate of the number of angles computed initially with respect to the total number of angles between all the views is 4.5%: a very low value, due to the setting being large and densely populated with reference images. However, this rate grows gradually to reach a high final value of 82.9%. In the end, only 20 out of the 213 available images could not be connected to the rest of the set, which means that 90.1% of the images were interconnected. These experimental results show the ability of our technique to connect robustly the views between different spaces (e.g. different rooms), through large numbers of intermediate images, and across long distances in a given environment. They also illustrate how our approach can perform well in scenarios with low quality images and few matched features.

3) *Outdoor experiment*: The method was also tested in an outdoor scenario, using a data set acquired in a parking lot at our university campus. The images were obtained using the same robotic platform and camera arrangement described in section IV-B.1. The separation between the locations where the different reference images were taken was much larger than in the indoors tests. For this experiment we used a set of 18 images acquired from positions forming a rectangular grid of approximate size  $62 \times 31 m^2$ . Thus, the distance between opposite corners of the grid was of around 70 m. The minimum distance between two images in this set was 12.5 m. Once again, the SIFT keypoint extractor was employed to obtain feature matches between the images. The number of three-view correspondences used to compute the 1D trifocal tensors was in this case in the range of 20 to

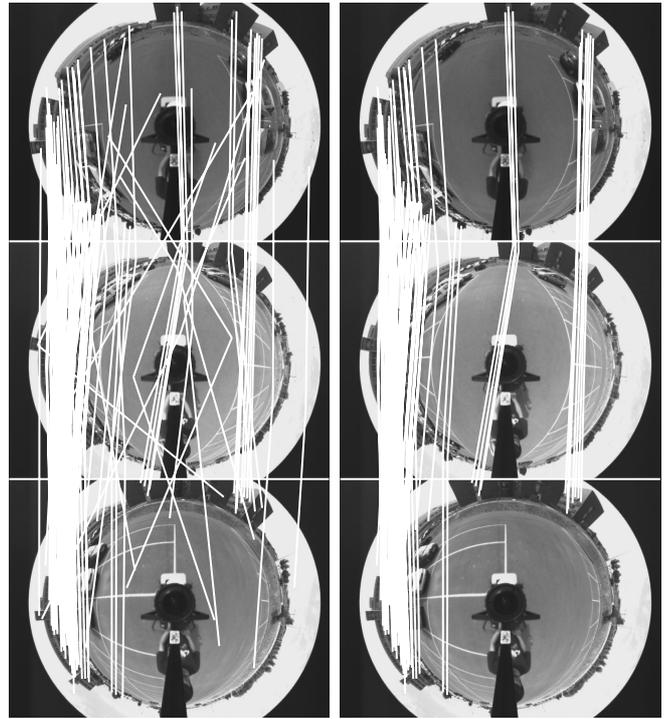


Fig. 14. Example of a trio of images from the outdoor data set with their putative SIFT correspondences joined by lines (left). Feature matches remaining after the robust computation of the 1D trifocal tensor (right).

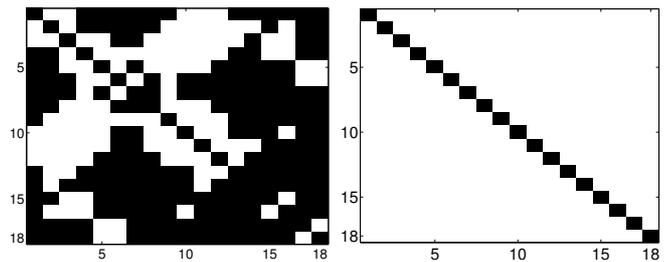


Fig. 15. Color-coded adjacency matrices for the outdoor experiment's image set. White color means '1', black color means '0'. The initial (left) and final (right) adjacency matrices are shown.

60. Useful point matches were found mainly on the walls of buildings, in the outer regions of the omnidirectional images. Figure 14 shows an example of the three-view matching process for this set, illustrating once again how the robust computation of the 1D trifocal tensor makes it possible to eliminate wrong correspondences.

Similarly to what occurred with the indoor tests, in this experiment only the connections (i.e. the angles of the epipoles) between physically close locations could be computed directly, with the remaining angles being calculated indirectly. We computed in this case 40% of the angles between the available views directly. The initial adjacency matrix obtained for this set of images is shown in Fig. 15, along with the resulting final adjacency matrix. As can be seen, in this case we eventually obtained a complete graph (100% of the angles, i.e. we were able to link all the reference positions between one another). The results of the angle

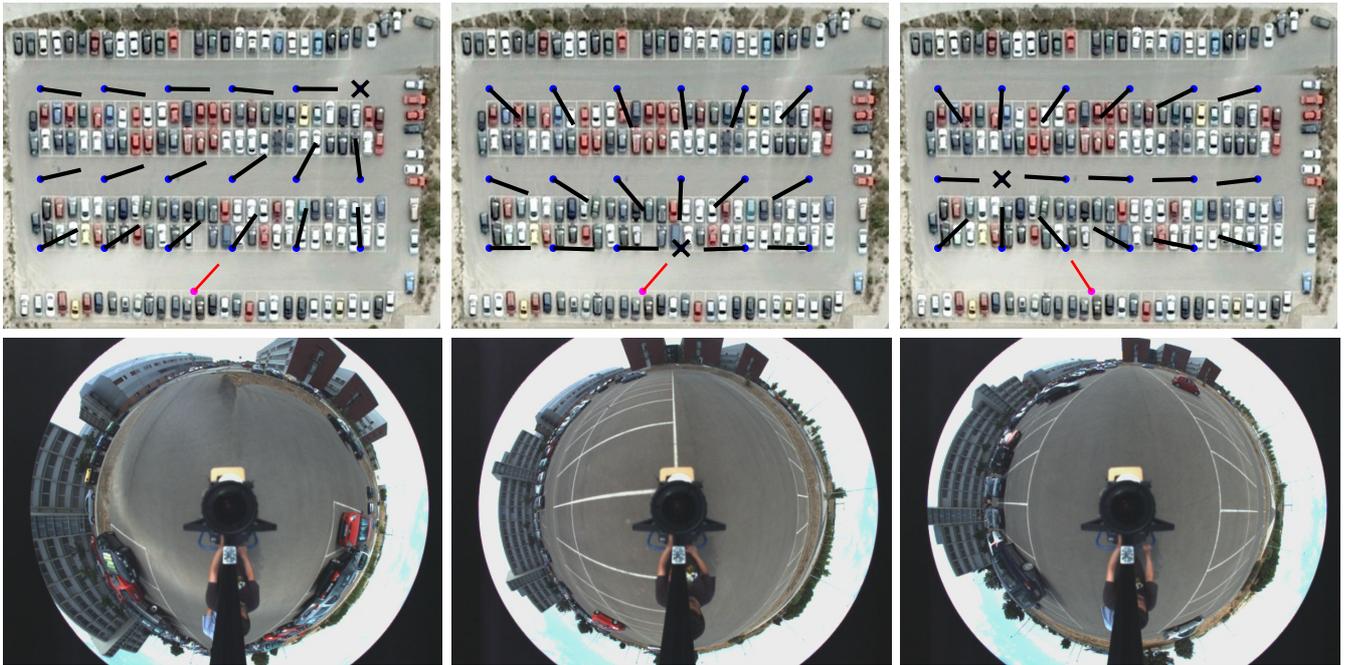


Fig. 16. Top: homing vectors estimated at every reference position for three different goal locations (marked with a cross) for the outdoor experiment. The homing vector from an outer position associated to an image not belonging to the reference grid is also shown. The vectors are plotted superimposed on a bird's-eye picture of the real setting. Bottom: goal images corresponding to the three homing vector representations on top.

calculations for the outdoor experiment are illustrated in Fig. 16, which shows the computed homing vectors from the reference positions to the goal (marked with a cross) for three different cases. The actual target images corresponding to each case are also shown in the figure. We also show the homing vectors from a position outside of the grid of reference views. This illustrates that, as long as it is possible to find sufficient three-view feature matches to compute the 1D trifocal tensor reliably between a given initial image and two of the reference images, our homing method will work, independently of the physical location where that initial image is captured. The setting appears full of cars in the bird's eye view on which we superimpose our results, but it was much less crowded when the images were acquired. We can see that the results are good in terms of robustness and accuracy. As previously commented in section IV-B.2, in order to navigate in this setting we would need to ensure that the goal location is reachable, e.g. by defining intermediate homing targets such that the robot always moves along the corridors of the parking lot.

The results of the outdoor experiment show that even with a relatively sparse representation of the environment, in terms of the separation between reference images, our approach is still capable of performing robustly. It has potential to cover and interconnect large areas, and allows homing to targets at long distances.

The different experiments we have conducted demonstrate that our method is versatile and has the ability to perform robustly in different scenarios, both indoors and outdoors.

## V. DISCUSSION

In this section we compare the characteristics of our homing method with those of other existing approaches, and discuss a number of aspects related to it.

### A. Comparison with existing work

In the following, we discuss the differences between the approach we present and other closely related works in the fields of visual homing and navigation. It is difficult to make a homogeneous and fair performance comparison among methods, since each of them has very different characteristics. Still, a number of aspects can be addressed. In order to allow a more specific and detailed comparative analysis, let us take at certain points of this discussion the method by Argyros et al. [10] as a representative of visual memory-based approaches, and the work by Goedemé et al. [27] as an example of topological map-based techniques.

The method [10] is, like ours, purely angle-based. The shape of the generated path to the goal depends heavily on the particular distribution of the scene landmarks. In contrast, in our method the path is independent of this. The difference in behavior is illustrated in the simulation we present in Fig. 17. Also, [10] can only be used to replay previously executed homing paths (i.e. it is a pure path-following approach). The method [27] uses an initial homing vector which is refined during the execution of the homing task. This initial vector is computed from the decomposition of the essential matrix for omnidirectional cameras, which is also used for the estimation of the homing direction in the work [28]. We provide in Fig. 18 a comparison of the error in the calculation of the homing vector when it is carried

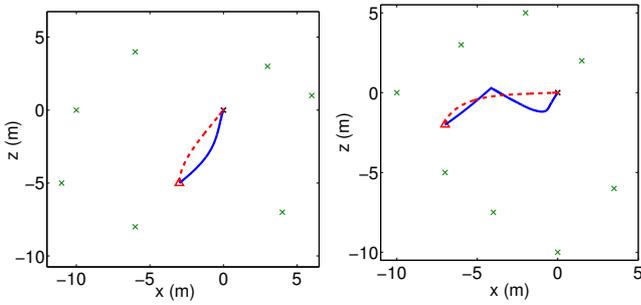


Fig. 17. Homing paths from two different initial locations to the goal location situated in  $(0,0)$ , with seven (left) and eight (right) feature matches. The paths generated by the method proposed by Argyros et al. [10] (solid line) and the paths generated by the approach presented in this paper (dashed line) are shown.

out using our method (based on the 1D trifocal tensor) and using the essential matrix. As can be seen, the trifocal tensor based-method is significantly more robust against noise.

As occurs with most existing visual homing approaches, both [10] and [27] use only two views to carry out the task, while we use the relations between three views. In these two approaches, the images employed can be acquired in any spatial distribution, whereas in our method, care must be taken to avoid acquiring all the images from positions on a straight line. Both of the related techniques take measures to prevent outlier matches, but these are based on two-view constraints, while we employ a stronger, three-view constraint (the 1D trifocal tensor). In contrast with our approach, none of the other two methods takes explicitly into account nonholonomic vehicle kinematics in the design of the motion strategy.

The approach we propose needs a minimum number of four reference images (the goal image and three others). It requires seven three-view correspondences in general configuration (i.e. not all on the same plane of the scene), so that the 1D trifocal tensor can be computed. This minimum number of correspondences is reduced to five if the camera is calibrated. Our second indoor experiment shows that our approach is able to operate with a minimum of twelve correct point correspondences between three views, although a number higher than twenty is desirable for better accuracy. In comparison, the method by Argyros et al. [10] employs a minimum of three matches between two views, although in practice they require several tens of correspondences to achieve good performance. The technique by Goedemé et al. [27] needs a minimum of eight points in general configuration matched between two views. Again, more points are usually considered, so as to make the computations more reliable. We believe the information requirements of our method, although higher than those of other existing approaches, are very reasonable, since image feature extractors typically provide many more matches than the minimum needed.

The fact that long-range homing can be carried out directly in our method, provides advantages with respect to both image memory-based and topological map-based long-range navigation approaches, since we will require fewer homing

sub-paths in general, thus increasing the efficiency. We believe that our method can also be more flexible if changes in the path are required: unlike in [10], [17], [30], the path to a given goal location is not restricted to following a fixed sequence of images, and unlike in [27], [28], all the positions of the reference views may be chosen as direct homing goals at all times.

We believe that two advantages of our method are its robustness, thanks to the use of the trifocal constraint, the geometric consistency checks and the stability properties of the control law, and its accuracy, due to the fact that both the employed visual information and the designed control law are purely angular.

### B. Practical and performance-related considerations

Next, we address several points related to the performance of our method and discuss a series of practical considerations. An interesting issue to tackle is how the error in the computed angles propagates as new angles are obtained using the indirect estimation procedure (section II-B). The theoretical analysis of this aspect in a general case (i.e. for arbitrary angles in the four-view set of Fig. 3) turns out to be very involved. We have observed through extensive experiments that the angular error does not propagate significantly across the graph we construct. Moreover, note that large error propagation is prevented by the three-view geometric consistency checks we perform, which ensure that the errors in the angle computations are always maintained within certain bounds. We provide as example the result from a simulation for 4, 8 and 16 images in a square grid-shaped reference set in Fig. 18. For this simulation, only the minimum necessary angles were computed directly, and all the rest were obtained indirectly. Notice that the average error does not increase when computing the indirect angles.

Regarding the performance of the angular disambiguation process described in section II-A.2, we have verified through extensive experiments that the procedure we propose has a very high success ratio. Figure 18 illustrates some of our results. As can be seen, and for obvious reasons, the disambiguation of the assignment of the epipoles has a larger probability of failure when the positions of the images are closer to being along a straight line. We show in the figure that the epipoles can frequently be assigned wrongly in quite extreme cases (very few points matched, and very high noise). With just a few more points (twenty) available, the success rate improves drastically. We have observed through experimentation that the ambiguity in determining the relative motion between two views from the 1D epipoles (Fig. 2) is resolved very effectively by our proposed method. Even for the minimum possible number of points (7) and very high noise ( $5^\circ$  of standard deviation in the projected angles), the disambiguation turns out to be correct in 99% of the cases. Note that the incorrect two-view disambiguations that may originate from errors in the assignment of the epipoles are robustly detected and discarded when the joint, three-view coherence of the two-view results is checked (algorithm 1).

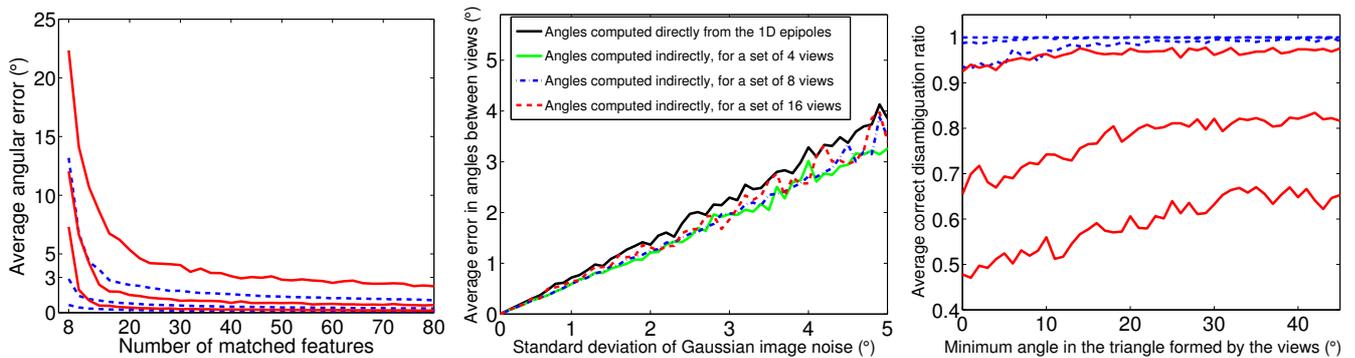


Fig. 18. Left: average error in the computed angle of the homing vector vs. number of matched features used. Solid line, top to bottom curves: homing vector computed from the essential matrix, with added Gaussian noise of standard deviation 3, 1 and  $0.3^\circ$ , respectively. Dashed line, top to bottom curves: homing vector computed through our method based on the 1D trifocal tensor, with added Gaussian noise of standard deviation 3, 1 and  $0.3^\circ$ , respectively. Center: Error propagation in the computed angles between views, for reference sets of 4, 8 and 16 images. Right: performance of the procedure to disambiguate the assignment of the epipoles. Solid line, top to bottom curves: success ratio of the disambiguation vs. degree of alignment of the locations of the views, for seven matched points and Gaussian noise of standard deviation 1, 3 and  $5^\circ$ , respectively. Dashed line, top to bottom curves: success ratio for twenty matched points and Gaussian noise of standard deviation 1, 3 and  $5^\circ$ , respectively.

The required density of reference images with our method is dictated by the need for a sufficient number of three-view feature matches in order to compute the 1D trifocal tensor reliably. Thus, the maximum separation between the locations of two reference images, or the maximum distance from a given starting location to the locations where the images in the set were taken, is given by the distance at which the number of three-view matches drops below a certain threshold. Also, the positions of the reference images in our method must not all be on a straight line, since both the disambiguation and indirect angle computation procedures rely on the differences between the angles in the triangles created by groups of three views. While acquiring the reference images, it is simple to take this point into account and guarantee that their locations are distributed in the environment with sufficient diversity.

With these conditions in mind, we believe that it would be possible to do the exploration of the environment in order to capture the reference images in a semi-automatic manner. The robot could be run across the environment capturing images either at a certain rate, or when the number of matches fell below a certain threshold. A sufficient number of three-view matches between views needs to be available, and in practice, the parameters of the feature extractor/matcher have to be tuned in order to adapt to the conditions of the environment and the images. These adjustments, which can be typically done automatically, allow to optimize the number of correspondences. If dynamic changes in the environment occur, our homing approach will be robust to them as long as sufficient features between images can be matched.

The definition of the intermediate homing objectives for our method could be done manually, or aided by the use of image classification and place segmentation techniques, in the manner of [35], [36]. In our case, the information that could be used for these purposes is the initial connectivity graph of the environment. For a given view, this graph determines robustly what other views share significant common

visual content with it. In order to implement our method, it is required to integrate the homing approach with a reactive obstacle avoidance procedure, as is common in other works in the literature [10], [27].

## VI. CONCLUSION AND FUTURE WORK

We have presented a visual homing method for a robot moving on the ground plane. Employing omnidirectional images acquired from the current robot position, the goal location and a set of other positions in the environment, the method works by computing the relative angles between all the locations by means of the 1D trifocal tensor. We have proposed a control law to drive the robot to the goal employing the calculated angles. This law has been proven to possess strong stability properties. In addition, the experiments have shown that the approach provides good accuracy and can perform robustly both indoors and outdoors, with low image quality and with varying density of images in the reference set. Homing between distant locations and different settings is also feasible. The online operation of the method is fast, the feature extraction process being the main limitation in terms of speed. The method can be directly applied in settings where stored image databases are available.

We can think of a number of directions for future work. First, the method could be extended by adapting the control law so that it applies to vehicles with different motion constraints, in particular with car-like kinematics. Also, it would be interesting to consider the possible methods to cluster the images in the constructed graph to allow the definition of the intermediate targets for a long-range homing task. Another possible issue to address would be how the acquisition of the reference image set could be carried out in an unknown environment, taking into account such aspects as the exploration strategy, the desired coverage area, the spatial distribution of images and the separation between them. Integrating the system with a reactive navigation method for obstacle avoidance would be another interesting problem. Finally, we believe it is feasible to extend the approach in

order to make it capable of driving the robot to positions in the environment defined not by images captured from them, but by visual features.

## APPENDIX

### A. Global asymptotic stability

*Theorem 1:* The system under the proposed control law (11), (12) is globally asymptotically stable if  $k_\omega > k_v \cdot \pi/d_{min}$ .

*Proof:* We will use Lyapunov techniques [37] to analyze the stability of the system. We define the following definite positive candidate Lyapunov function:

$$V(\mathbf{x}, t) = \frac{\rho^2}{2} + \frac{(\alpha_{CG} - \alpha_{CG}^d)^2}{2}, \quad (13)$$

where  $\rho$  is the distance between the current and goal positions, and  $\mathbf{x}$  is the state of the system, determined by  $\rho$  and  $\alpha_{CG}$ . The two state variables we use are a suitable choice, since we are only interested in reaching the goal position, regardless of the final orientation of the robot. As can be seen, both  $V$  and  $\dot{V}$  are continuous functions.

We note at this moment that the equilibrium in our system occurs at the two following points:  $(\rho, \alpha_{CG}) = (0, 0)$  and  $(\rho, \alpha_{CG}) = (0, \pi)$ , which correspond to the situations where the robot reaches the goal moving forwards or backwards, respectively. In order to account for the multiple equilibria, in the following we use the global invariant set theorem [38] to prove the asymptotic stability of the system.

What we need to show is that  $V$  is radially unbounded and  $\dot{V}$  is negative semi-definite over the whole state space. It is straightforward that  $V(\mathbf{x})$  is radially unbounded, given that  $V(\mathbf{x}) \rightarrow \infty$  as  $\|\mathbf{x}\| \rightarrow \infty$ . Next, we prove that the derivative  $\dot{V}(\mathbf{x})$  is negative definite. For our chosen candidate Lyapunov function, this derivative is as follows:

$$\dot{V} = \rho \dot{\rho} + (\alpha_{CG} - \alpha_{CG}^d) \dot{\alpha}_{CG}. \quad (14)$$

We will suppose that the vehicle on which the control method is to be implemented is a nonholonomic unicycle platform. The dynamics of the system as a function of the input velocities is then given, using the derivatives in polar coordinates with the origin at the goal, by  $\dot{\rho} = -v \cos(\alpha_{CG})$  and  $\dot{\alpha}_{CG} = -\omega + v \sin(\alpha_{CG})/\rho$ . Using the control velocities (11), (12) we obtain

$$\begin{aligned} \dot{V} = & -k_v \rho \cos^2(\alpha_{CG}) \frac{1}{n} \sum_{i=1}^n S_i - k_\omega (\alpha_{CG} - \alpha_{CG}^d)^2 \\ & + (\alpha_{CG} - \alpha_{CG}^d) \cdot \sin \alpha_{CG} \frac{k_v}{\rho} \cos \alpha_{CG} \cdot \frac{1}{n} \sum_{i=1}^n S_i. \end{aligned} \quad (15)$$

By definition  $\rho \geq 0$  and  $S_i \geq 0$ . It is then straightforward to see that the first and the second term of (15) are negative definite. However, the third term can be positive. The interpretation is that for the system to be stable, the convergence speed provided by the angular velocity has to be higher than the convergence speed given by the linear

velocity. Otherwise, the angular error is not corrected fast enough and the robot will move following spirals around the goal. Still, the stability can be guaranteed if the control gains are selected properly. From (15) we can see that it is guaranteed that  $\dot{V} < 0$  if the following inequality holds:

$$|k_\omega \cdot (\alpha_{CG} - \alpha_{CG}^d)| > |\sin(\alpha_{CG}) \cos(\alpha_{CG}) \frac{k_v}{\rho} \cdot \frac{1}{n} \sum_{i=1}^n S_i|. \quad (16)$$

This is equivalent to the following condition on the angular velocity gain:

$$k_\omega > \left| \frac{\sin(\alpha_{CG})}{(\alpha_{CG} - \alpha_{CG}^d)} \cos(\alpha_{CG}) \cdot k_v \cdot \frac{1}{n} \sum_{i=1}^n \frac{S_i}{\rho} \right|. \quad (17)$$

We aim to find an upper bound to the right side of (17). We start by analyzing the first fraction. Since  $\alpha_{CG}^d$  is equal to either 0 or  $\pi$ , and  $\sin(\alpha_{CG}) = -\sin(\alpha_{CG} - \pi)$ , we have:

$$\left| \frac{\sin(\alpha_{CG})}{(\alpha_{CG} - \alpha_{CG}^d)} \right| = \left| \frac{\sin(\alpha_{CG})}{(\alpha_{CG})} \right| \leq 1, \quad (18)$$

as  $\sin(\alpha_{CG})/\alpha_{CG}$  is a *sinc* function, whose maximum absolute value occurs at  $\alpha_{CG} = 0$  and equals 1. We now look for a bound to the  $S_i/\rho$  term in (17). The angular sector  $S_i$  seen from reference view  $i$  has a value lying in the interval  $0 \leq S_i \leq \pi$ . We will study two subintervals separately:

- $0 \leq S_i \leq \pi/2$ . Applying the law of sines on the triangle defined by vertices  $C$ ,  $G$  and  $R_i$  in Fig. 5, the addend in (17) corresponding to reference view  $i$  becomes:

$$\frac{S_i}{\rho} = \frac{S_i}{\sin(S_i)} \cdot \frac{\sin(\widehat{CR_iG})}{d_i} \leq \frac{\pi}{2 \cdot d_{min}} \quad (19)$$

The first fraction of the product in (19) is a function of  $S_i$  whose value equals 1 at  $S_i = 0$  and increases monotonically to a value of  $\pi/2$  at  $S_i = \pi/2$ , which is the limit of the interval we are considering. Since the second fraction has an upper bound equal to  $1/d_{min}$ , the product of the two is upper-bounded by  $\pi/(2 \cdot d_{min})$ .

- $\pi/2 < S_i \leq \pi$ . In this case,  $\rho > d_i$ , and an upper bound is readily found for the addend in (17) corresponding to reference view  $i$ :

$$\frac{S_i}{\rho} \leq \frac{\pi}{d_{min}}. \quad (20)$$

Thus, the contribution of each of the reference views to the sum is upper-bounded by the higher of the two bounds in (19) and (20), which is  $\frac{\pi}{d_{min}}$ . The mean of all the individual contributions is therefore bounded by this value, i.e.:

$$\frac{1}{n} \sum_{i=1}^n \frac{S_i}{\rho} \leq \frac{\pi}{d_{min}}, \quad (21)$$

and inequality (17) becomes:

$$k_\omega > \frac{k_v \cdot \pi}{d_{min}}. \quad (22)$$

■

## B. Local exponential stability

*Proposition 1:* The system under the proposed control law (11), (12) is locally exponentially stable.

*Proof:*

We analyze the behavior of the system locally, i.e. assuming the orientation of the robot has already been corrected ( $\alpha_{CG} = \alpha_{CG}^d$ ). The dynamics of the distance from the goal for the unicycle vehicle considered is then given by:

$$\dot{\rho} = -v \cos \alpha_{CG} = -k_v \frac{1}{n} \sum_{i=1}^n S_i. \quad (23)$$

Now, taking into account that  $S_i \geq \sin S_i$  in all the interval of possible values ( $0 \leq S_i \leq \pi$ ), we have:

$$\dot{\rho} \leq \frac{-k_v}{n} \sum_{i=1}^n \sin S_i = - \left[ \frac{k_v}{n} \sum_{i=1}^n \frac{\sin(\widehat{CR_iG})}{d_i} \right] \cdot \rho. \quad (24)$$

It can be readily seen, looking at Fig. 5, that for any given current position  $C$  of the robot,  $\sin(\widehat{CR_iG})$  will be greater than zero for at least one  $R_i$  as long as there are at least three reference views (including the goal) and their locations are not collinear. Thus, there exists a positive value  $\lambda_{min}$  such that

$$\frac{k_v}{n} \sum_{i=1}^n \frac{\sin(\widehat{CR_iG})}{d_i} \geq \lambda_{min} > 0. \quad (25)$$

From (24) and (25) it can be concluded that the local convergence to the target state is bounded by an exponential decay, i.e. the system is locally exponentially stable. ■

## REFERENCES

- [1] G. N. DeSouza and A. C. Kak, "Vision for mobile robot navigation: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 2, pp. 237–267, 2002.
- [2] F. Chaumette and S. Hutchinson, "Visual servo control, part I: Basic approaches," *IEEE Robotics and Automation Magazine*, vol. 13, no. 4, pp. 82–90, 2006.
- [3] G. Chesi and K. Hashimoto, Eds., *Lecture notes in control and information sciences: Vol. 401*. Springer, 2010, ch. Visual Servoing via Advanced Numerical Methods.
- [4] D. Lambrinos, R. Möller, T. Labhart, R. Pfeifer, and R. Wehner, "A mobile robot employing insect strategies for navigation," *Robotics and Autonomous Systems*, vol. 30, no. 1-2, pp. 39–64, 2000.
- [5] K. Weber, S. Venkatesh, and M. V. Srinivasany, "Insect-inspired robotic homing," *Adaptive Behavior*, vol. 7, no. 1, pp. 65–97, 1998.
- [6] R. Möller, A. Vardy, S. Krefit, and S. Ruwisch, "Visual homing in environments with anisotropic landmark distribution," *Autonomous Robots*, vol. 23, no. 3, pp. 231–245, 2007.
- [7] D. Churchill and A. Vardy, "Homing in scale space," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2008, pp. 1307–1312.
- [8] W. Stürzl and H. A. Mallot, "Efficient visual homing based on Fourier transformed panoramic images," *Robotics and Autonomous Systems*, vol. 54, no. 4, pp. 300–313, 2006.
- [9] J. Hong, X. Tan, B. Pinette, R. Weiss, and E. M. Riseman, "Image-based homing," *Control Systems Magazine, IEEE*, vol. 12, no. 1, pp. 38–45, 1992.
- [10] A. A. Argyros, K. E. Bekris, S. C. Orphanoudakis, and L. E. Kavradi, "Robot homing by exploiting panoramic vision," *Autonomous Robots*, vol. 19, no. 1, pp. 7–25, 2005.
- [11] J. Lim and N. Barnes, "Robust visual homing with landmark angles," in *Proceedings of Robotics: Science and Systems*, 2009.
- [12] R. Basri, E. Rivlin, and I. Shimshoni, "Visual homing: Surfing on the epipoles," *International Journal of Computer Vision*, vol. 33, no. 2, pp. 117–137, 1999.
- [13] G. Chesi and K. Hashimoto, "A simple technique for improving camera displacement estimation in eye-in-hand visual servoing," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 26, no. 9, pp. 1239–1242, 2004.
- [14] G. López-Nicolás, C. Sagüés, J. Guerrero, D. Kragic, and P. Jensfelt, "Switching visual control based on epipoles for mobile robots," *Robotics and Autonomous Systems*, vol. 56, no. 7, pp. 592–603, 2008.
- [15] G. López-Nicolás and C. Sagüés, "Vision-based exponential stabilization of mobile robots," *Autonomous Robots*, vol. 30, no. 3, pp. 293–306, 2011.
- [16] J. Chen, W. Dixon, M. Dawson, and M. McIntyre, "Homography-based visual servo tracking control of a wheeled mobile robot," *IEEE Transactions on Robotics*, vol. 22, no. 2, pp. 407–416, 2006.
- [17] J. Courbon, Y. Mezouar, and P. Martinet, "Indoor navigation of a non-holonomic mobile robot using a visual memory," *Autonomous Robots*, vol. 25, no. 3, pp. 253–266, 2008.
- [18] G. López-Nicolás, J. J. Guerrero, and C. Sagüés, "Multiple homographies with omnidirectional vision for robot homing," *Robotics and Autonomous Systems*, vol. 58, no. 6, pp. 773–783, 2010.
- [19] R. I. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*, 2nd ed. Cambridge University Press, 2004.
- [20] G. López-Nicolás, J. J. Guerrero, and C. Sagüés, "Visual control through the trifocal tensor for nonholonomic robots," *Robotics and Autonomous Systems*, vol. 58, no. 2, pp. 216–226, 2010.
- [21] A. Shademan and M. Jägersand, "Three-view uncalibrated visual servoing," in *IEEE International Conference on Intelligent Robots and Systems*, 2010, pp. 6234–6239.
- [22] L. Quan, "Two-way ambiguity in 2D projective reconstruction from three uncalibrated 1D images," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 23, no. 2, pp. 212–216, 2001.
- [23] F. Dellaert and A. W. Stroupe, "Linear 2D localization and mapping for single and multiple robot scenarios," in *IEEE International Conference on Robotics and Automation*, 2002, pp. 688–694.
- [24] J. J. Guerrero, A. C. Murillo, and C. Sagüés, "Localization and matching using the planar trifocal tensor with bearing-only data," *IEEE Transactions on Robotics*, vol. 24, no. 2, pp. 494–501, 2008.
- [25] M. Aranda, G. López-Nicolás, and C. Sagüés, "Omnidirectional visual homing using the 1D trifocal tensor," in *IEEE International Conference on Robotics and Automation*, 2010, pp. 2444–2450.
- [26] H. Becerra, G. Lopez-Nicolas, and C. Sagüés, "Omnidirectional visual control of mobile robots based on the 1D trifocal tensor," *Robotics and Autonomous Systems*, vol. 58, no. 6, pp. 796–808, 2010.
- [27] T. Goedemé, M. Nuttin, T. Tuytelaars, and L. Van Gool, "Omnidirectional vision based topological navigation," *International Journal of Computer Vision*, vol. 74, no. 3, pp. 219–236, 2007.
- [28] O. Booij, B. Terwijn, Z. Zivkovic, and B. Kröse, "Navigation using an appearance based topological map," in *IEEE International Conference on Robotics and Automation*, 2007, pp. 3927–3932.
- [29] M. O. Franz, B. Schölkopf, P. Georg, H. A. Mallot, and H. H. Bülthoff, "Learning view graphs for robot navigation," *Autonomous Robots*, vol. 5, no. 1, pp. 111–125, 1998.
- [30] A. Cherubini and F. Chaumette, "Visual Navigation With Obstacle Avoidance," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2011, pp. 1593–1598.
- [31] K. Åström and M. Oskarsson, "Solutions and ambiguities of the structure and motion problem for 1D retinal vision," *Journal of Mathematical Imaging and Vision*, vol. 12, no. 2, pp. 121–135, 2000.
- [32] D. Lowe, "Distinctive image features from scale-invariant keypoints," *International Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [33] A. Shashua and M. Werman, "Trilinearity of three perspective views and its associated tensor," in *International Conference on Computer Vision*, 1995, pp. 920–925.
- [34] O. Booij, Z. Zivkovic, and B. Kröse, "Sparse appearance based modeling for robot localization," in *IEEE International Conference on Intelligent Robots and Systems*, 2006, pp. 1510–1515.
- [35] Z. Zivkovic, O. Booij, and B. Kröse, "From images to rooms," *Robotics and Autonomous Systems*, vol. 55, no. 5, pp. 411–418, 2007.
- [36] Z. Zivkovic, B. Bakker, and B. Kröse, "Hierarchical map building and planning based on graph partitioning," in *IEEE International Conference on Robotics and Automation*, 2006, pp. 803–809.
- [37] H. K. Khalil, *Nonlinear Systems*, 3rd ed. Prentice Hall, 2001.
- [38] J.-J. E. Slotine and W. Li, *Applied Nonlinear Control*. Prentice Hall, 1991.