

# Monocular Omnidirectional Visual Odometry for Outdoor Ground Vehicles

Davide Scaramuzza and Roland Siegwart

Swiss Federal Institute of Technology Zurich (ETHZ)  
Autonomous System Laboratory (ASL)  
CH-8092, Zurich, Switzerland  
{davide.scaramuzza,r.siegwart}@ieee.org  
<http://www.asl.ethz.ch>

**Abstract.** This paper describes an algorithm for visually computing the ego-motion of a vehicle relative to the road under the assumption of planar motion. The algorithm uses only images taken by a single omnidirectional camera mounted on the roof of the vehicle. The front ends of the system are two different trackers. The first one is a homography-based tracker that detects and matches robust scale invariant features that most likely belong to the ground plane. The second one uses an appearance based approach and gives high resolution estimates of the rotation of the vehicle. This 2D pose estimation method has been successfully applied to videos from an automotive platform. We give an example of camera trajectory estimated purely from omnidirectional images over a distance of 400 meters. For performance evaluation, the estimated path is superimposed onto an aerial image. In the end, we use image mosaicing to obtain a textured 2D reconstruction of the estimated path.

**Keywords:** omnidirectional camera, visual odometry, vehicle ego-motion estimation, homography, SIFT features.

## 1 Introduction

Accurate estimation of the ego-motion of a vehicle relative to the road is a key component for autonomous driving and computer vision based driving assistance. Most of the work in estimating robot motion has been produced using stereo cameras and can be traced back to Moravec's work [1]. Similar work has been reported also elsewhere (see [2, 3]). Furthermore, stereo visual odometry has also been successfully used on Mars by the NASA rovers since early 2004 [4]. Nevertheless, visual odometry methods for outdoor applications have been also produced, which use a single camera alone. Very successful results have been obtained over long distances using either perspective or omnidirectional cameras (see [3, 5]). In [3], the authors deal with the case of a stereo camera but they also provide a monocular solution implementing a structure from motion algorithm that takes advantage of the 5-point algorithm and RANSAC robust estimation [14]. In [5], the authors provide two approaches for monocular visual odometry

based on omnidirectional imagery. In the first approach, they use optical flow computation while in the second one structure from motion.

In our approach, we use a single calibrated omnidirectional camera mounted on the roof of the car (Fig. 3). We assume that the vehicle undergoes a purely two-dimensional motion over a predominant flat ground. Furthermore, because we want to perform visual odometry in city streets, flat terrains, as in well as in motorways where buildings or 3D structure are not always present, we estimate the motion of the vehicle by tracking the ground plane.

Ground plane tracking has been already exploited by the robotics community for indoor visual navigation and most works have been produced using standard perspective cameras ([6, 7, 8, 9]). In those works, the motion of the vehicle is estimated by using the property that the projection of the ground plane into two different camera views is related by a homography.

In this paper, we propose a similar approach for central omnidirectional cameras but our goal is to estimate the ego-motion of the vehicle in outdoor environments and over long distances. Thanks to the large field of view of the panoramic camera, SIFT keypoints [10] from the scene all around the car are extracted and matched between consecutive frames. After RANSAC based outlier removal [14], we use these features only to compute the translation in the heading direction. Conversely, to estimate the rotation angle of the vehicle we use an appearance based method. We show that by adding this second approach the result outperforms the pure feature based approach.

This paper is organized as follows. Section 2 describes our homography based ground plane navigation. Section 3 describes the appearance based method. Section 4 details the steps of the whole visual odometry algorithm. Finally, Section 5 is dedicated to the experimental results.

## 2 Homography Based Ground Plane Navigation

The motion information that can be extracted by tracking 2D features is central to our vehicle navigation system. Therefore, we briefly review here a method that uses planar constraints and point tracking to compute the motion parameters.

### 2.1 Homography and Planar Motion Parameters

Early work on exploiting coplanar relations has been presented by Tsai and Huang [11], Longuet-Higgins [12], and Faugeras and Lustman [13]. The coplanar relation between two different views of the same plane can be summarized as follows.

$$\lambda \mathbf{x}_2 = \mathbf{K}(\mathbf{R} + \frac{\mathbf{T}\mathbf{n}^T}{h})\mathbf{K}^{-1}\mathbf{x}_1 = \mathbf{H}\mathbf{x}_1 \quad (1)$$

where  $\mathbf{R} \in SO(3)$  and  $\mathbf{T} \in \mathbb{R}^3$  are the rotation and the translation matrices encoding the relative position of the two views;  $\mathbf{n} \in \mathbb{R}^3$  is the plane normal and  $h \in \mathbb{R}$  is the distance to the plane;  $\mathbf{x}_1, \mathbf{x}_2$  are the images of the same

scene points expressed in homogeneous coordinates  $([x, y, 1]^T)$ ;  $\mathbf{K}$  is a  $3 \times 3$  matrix describing the camera intrinsic parameters;  $\lambda$  is a scalar;  $\mathbf{H}$  is a  $3 \times 3$  matrix called homography that relates the two camera projections of the same plane points. Note that matrix  $\mathbf{K}$  in equation (1) is defined only for perspective cameras. However, in this paper we assume that our omnidirectional camera is already intrinsically calibrated and that the image points  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are already normalized to have the third component equal to 1. This allows us to write  $\mathbf{K} = \mathbf{I}$ . If not stated otherwise, in the remainder of this paper we will assume that the image coordinates are always normalized. To calibrate our omnidirectional camera, we used the method proposed in [15].

In our experiments, we mounted the omnidirectional camera on the roof of the car (as in Fig. 3) with the  $z$ -axis of the mirror perpendicular to the ground plane (Fig. 4). By fixing the origin of our coordinate system in the center of projection of the omnidirectional camera (Fig. 4), we have that  $\mathbf{n} = [0, 0, -1]^T$ . The distance  $h$  of the origin to the ground plane can be manually measured.

## 2.2 Decomposing $\mathbf{H}$

If a camera is internally calibrated, it is possible to recover  $\mathbf{R}$ ,  $\mathbf{T}$ , and  $\mathbf{n}$  from  $\mathbf{H}$  up to at most a two-fold ambiguity. A linear method for decomposing  $\mathbf{H}$  was originally developed by Wunderlich [17] and later reformulated by Triggs [18]. The algorithm of Triggs is based on the singular value decomposition of  $\mathbf{H}$ . The description of this method as well as its Matlab implementation can be found in [18]. This algorithm outputs two possible solutions for  $\mathbf{R}$ ,  $\mathbf{T}$ , and  $\mathbf{n}$  which are all internally self-consistent. In the general case, some false solutions can be eliminated by sign (visibility) tests or geometric constraints, while in our case we can disambiguate the solutions by choosing the one for which the computed plane normal  $\mathbf{n}$  is closer to  $[0, 0, -1]^T$ . In the remainder of this paper, we will refer to this method as the ‘‘Triggs algorithm’’.

## 2.3 Non-linear Refinement

The solution given by the Triggs algorithm is obtained by a linear method that minimizes an algebraic distance which is not physically meaningful. We can refine it through maximum likelihood inference. The maximum likelihood estimate can be obtained by minimizing the following functional:

$$\min_{\theta, t_1, t_2} \sum_{i=1}^n \|\mathbf{x}_1^i - \hat{\mathbf{x}}_1^i(\mathbf{R}, \mathbf{T}, \mathbf{n})\|^2 + \|\mathbf{x}_2^i - \hat{\mathbf{x}}_2^i(\mathbf{R}, \mathbf{T}, \mathbf{n})\|^2, \quad (2)$$

with  $\hat{\mathbf{x}}_1 = \mathbf{H}^{-1}\mathbf{x}_2$  and  $\hat{\mathbf{x}}_2 = \mathbf{H}\mathbf{x}_1$  according to equation (1). To minimize (2), we used the Levenberg-Marquadt algorithm. Furthermore, because we assume planar motion, we constraint the minimization so that the rotation is about the plane normal and the translation is parallel to the same plane.

## 2.4 Coplanarity Check

So far we have assumed that the corresponding image pairs  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are correctly matched and that their correspondent scene points are coplanar. Even though in omnidirectional images taken from the roof of the car the ground plane is predominant, there are also many feature points that come from other objects than just the road, like cars, buildings, trees, guardrails, etc. Furthermore, there are also many unavoidable false matches that are more numerous than those usually output by SIFT on standard perspective images (about 20-30% according to [10]) because of the large distortion introduced by the mirror. To discard the outliers, we use the Random Sample Consensus paradigm (RANSAC) [14].

## 3 Visual Compass

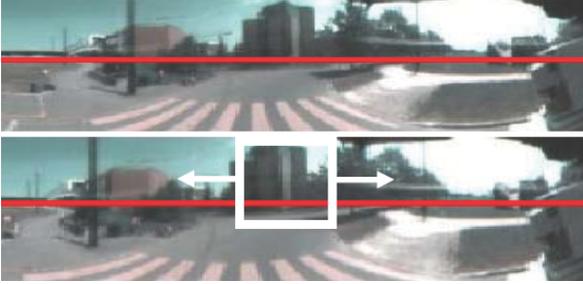
In the previous session, we described how to use point features to compute the rotation and translation matrices. Unfortunately, when using features to estimate the motion, the resulting rotation is extremely sensitive to systematic errors due to the intrinsic calibration of the camera or the extrinsic calibration between the camera and the ground plane. This effect is even more accentuated with omnidirectional cameras due to the large distortion introduced by the mirror. In addition to this, integrating rotational information over the time has the major drawback of generally becoming less and less accurate as integration introduces additive errors at each step. An example of camera trajectory recovered using only the feature based approach described in Section 2 is depicted in Fig. 4.

To improve the accuracy of the rotation estimation, we use an appearance based approach. This approach was inspired by the work of Labrosse [16], which describes a method to use omnidirectional cameras as visual compass.

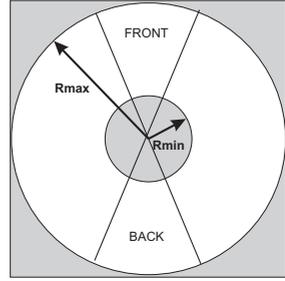
Directly using the appearance of the world as opposed to extracting features or structure of the world is attractive because methods can be devised that do not need precise calibration steps. Here, we describe how we implemented our visual compass.

For ease of processing, every omnidirectional image is unwrapped into cylindrical panoramas (Fig. 1). The unwrapping considers only the white region of the omnidirectional image that is depicted in Fig 2. We call these unwrapped versions “appearances”. If the camera is perfectly vertical to the ground, then a pure rotation about its vertical axis will result in a simple column-wise shift of the appearance in the opposite direction. The exact rotation angle could then be retrieved by simply finding the best match between a reference image (before rotation) and a column-wise shift of the successive image (after rotation). The best shift is directly related to the rotation angle undertaken by the camera. In the general motion, translational information is also present. This general case will be discussed later.

The input to our rotation estimation scheme is thus made of appearances that need to be compared. To compare them, we use the Euclidean distance. The



**Fig. 1.** Two unwrapped omnidirectional images taken at consecutive time stamps. For reasons of space, here only one half of the whole  $360\text{ deg}$  is shown. The red line indicates the horizon.



**Fig. 2.** The cylindrical panorama is obtained by unwrapping the white region

Euclidean distance between two appearances  $I_i$  and  $I_j$ , with  $I_j$  being column-wise shifted (with column wrapping) by  $\alpha$  pixels, is:

$$d(I_i, I_j, \alpha) = \sqrt{\sum_{k=1}^h \sum_{h=1}^w \sum_{l=1}^c |I_i(k, h, l) - I_j(k, h - \alpha, l)|^2} \quad (3)$$

where  $h \times w$  is the image size, and  $c$  is the number of color components. In our experiments, we used the RGB color space, thus having three color components per pixel.

If  $\alpha_m$  is the best shift that minimizes the distance between two appearances  $I_i$  and  $I_j$ , the rotation angle  $\Delta\vartheta$  (in degrees) between  $I_i$  and  $I_j$  can be computed as:

$$\Delta\vartheta = \alpha_m \cdot \frac{360}{w} \quad (4)$$

The width  $w$  of the appearance is the width of the omnidirectional image after unwrapping and can be chosen arbitrarily. In our experiments, we used  $w = 360$ , that means the angular resolution was 1 pixel per degree. To increase the resolution to  $0.1\text{ deg}$ , we used cubic spline interpolation with 0.1 pixel precision. We also tried larger image widths but we did not get any remarkable improvement in the final result. Thus, we used  $w = 360$  as the unwrapping can be done in a negligible amount of time.

The distance minimization in (3) makes sense only when the camera undergoes a pure rotation about its vertical axis, as a rotation corresponds to a horizontal shift in the appearance. In the real case, the vehicle is moving and translational component is present. However, the “pure rotation” assumption still holds if the camera undergoes small displacements or the distance to the objects (buildings, tree, etc.) is big compared to the displacement. In the other cases, this assumption does not hold for the whole image but an improvement that can be done over the theoretical method is to only consider parts of the images, namely the front and back part (Fig. 2). Indeed, the contribution to the optical flow

by the motion of the camera is not homogeneous in omnidirectional images; a forward/backward translation mostly contributes in the regions corresponding to the sides of the camera and very little in the parts corresponding to the front and back of the camera, while the rotation contributes equally everywhere.

Because we are interested in extracting the rotation information, only considering the regions of the images corresponding to the front and back of the camera allows us to reduce most of the problems introduced by the translation, in particular sudden changes in appearance (parallax).

According to the last considerations, in our experiments we use a reduced Field Of View (FOV) around the front and back of the camera (Fig. 2). A reduced field of view of about  $30\text{ deg}$  around the front part is shown by the white window in Fig. 1. Observe that, besides reducing the FOV of the camera in the horizontal plane, we operate a reduction of the FOV also in the vertical plane, in particular under the horizon line. The objective is to reduce the influence of the changes in appearance of the road. The resulting vertical FOV is  $50\text{ deg}$  above and  $10\text{ deg}$  below the horizon line (the horizon line is indicated in red in Fig. 1).

## 4 Motion Estimation Algorithm

As we already mentioned, the appearance based approach provides rotation angle estimates that are more reliable and stable than those output by the pure feature based approach. Here, we describe how we combined the rotation angle estimates of Section 3 with the camera translation estimates of Section 2.

In our experiments, the speed of the vehicle ranged between 10 and 20 Km/h while the images were constantly captured at 10 Hz. This means that the distance covered between two consecutive frames ranged between 0.3 and 0.6 meters. For this short distance, the kinematic model of the camera configuration  $(x, y, \theta)$ , which contains its 2D position  $(x, y)$  and orientation  $\theta$ , can be approximated in this way:

$$\begin{cases} x_{i+1} = x_i + \delta\rho_i \cos\theta \\ y_{i+1} = y_i + \delta\rho_i \sin\theta \\ \theta_{i+1} = \theta_i + \delta\theta_i \end{cases} \quad (5)$$

where we use  $\delta\rho = |\mathbf{T}| h$  and  $\delta\theta = \Delta\vartheta$ .  $|\mathbf{T}|$  is the length of the translation vector;  $h$  is the scale factor (i.e. in our case this is the height of the camera to the ground plane). The camera rotation angle  $\Delta\vartheta$  is computed as in (4). Observe that we did not use at all the rotation estimates provided by the feature based method of Section 2.

Now, let us resume the steps of our motion estimation scheme, which have been detailed in Section 2 and 3. Our omnidirectional visual odometry operates as follows:

1. Acquire two consecutive frames. Consider only the region of the omnidirectional image, which is between  $Rmin$  and  $Rmax$  (Fig. 2).
2. Extract and match SIFT features between the two frames. Use the double consistency check to reduce the number of outliers. Then, use the calibrated camera model to normalize the feature coordinates.

3. Use RANSAC to reject points that are not coplanar (Section 2.4).
4. Apply the Triggs algorithm followed by non-linear refinement described in Section 2 to estimate  $\mathbf{R}$  and  $\mathbf{T}$  from the remaining inliers.
5. Unwrap the two images and compare them using the appearance method described in Section 3. In particular, minimize (3), with reduced field of view, to compute the column-wise shift  $\alpha_m$  between the appearances and use (4) to compute the rotation angle  $\Delta\vartheta$ .
6. Use  $\delta\rho = |\mathbf{T}| h$  and  $\delta\theta = \Delta\vartheta$  and integrate the motion using (5).
7. Repeat from step 1.

## 5 Results

The approach proposed in this paper has been successfully tested on a real vehicle equipped with a central omnidirectional camera. A picture of our vehicle (a Smart) is shown in Fig 3.

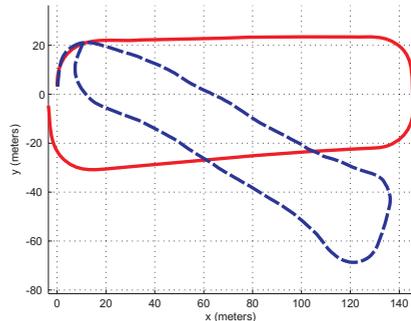
Our omnidirectional camera, composed of a hyperbolic mirror (KAIDAN 360 One VR) and a digital color camera (SONY XCD-SX910, image size  $640 \times 480$  pixels), was installed on the front part of the roof of the vehicle. The frames were grabbed at 10 Hz and the vehicle speed ranged between 10 and 20 Km/h.

The resulting path estimated by our visual odometry algorithm using a horizontal reduced FOV of 10 *deg* is shown in figures 4, 5, and 6. Our ground truth is a aerial image of the same test environment provided by Google Earth (Fig. 5). The units used in the three figures are in meters.

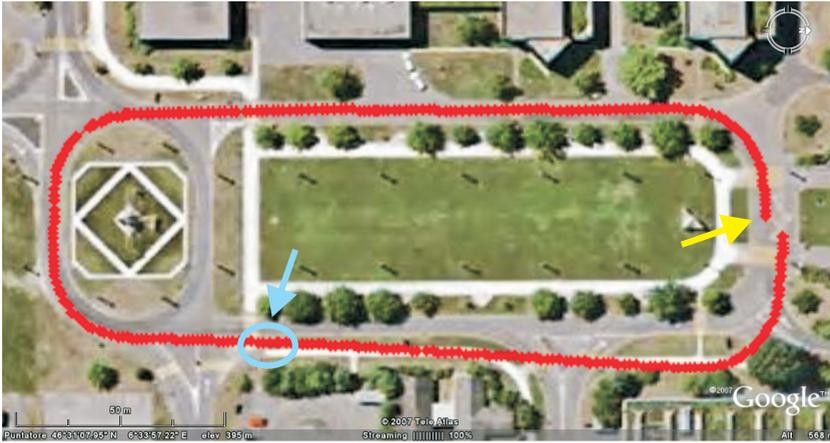
In this experiment, the vehicle was driven along a 400 meter long loop and returned to its starting position (pointed to by the yellow arrow in Fig. 5). The estimated path is indicated with red dots in Fig. 5 and is shown superimposed on the aerial image for comparison. The final error at the loop closure is about 6.5 meters. This error is due to the unavoidable visual odometry drift; however, observe that the trajectory is very well estimated until the third 90-degree turn.



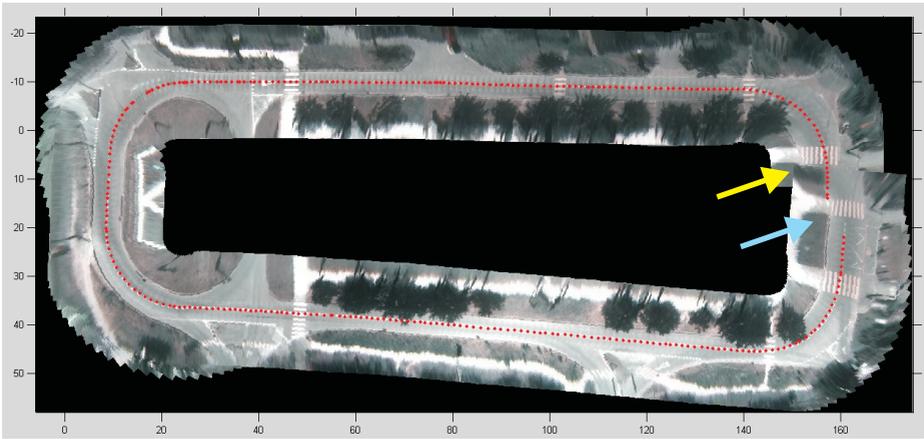
**Fig. 3.** Our vehicle with the omnidirectional camera (blue circle). The field of view is indicated by the red lines.



**Fig. 4.** Comparison between the standard feature based approach (dashed blue) and the approach combining features with visual compass proposed in this paper (red)



**Fig. 5.** The estimated path superimposed onto a Google Earth image of the test environment. The scale is shown at the lower left corner.



**Fig. 6.** Image mosaicing that shows a textured 2D reconstruction of the estimated path. The two arrows point out the final error at the loop closure (the pedestrian crossing pointed to by the cyan arrow should theoretically coincide with that pointed to by the yellow arrow).

After the third turn, the estimated path deviates smoothly from the expected path instead of continuing straight. After road inspection, we found that this deviation was due to three 0.3 meter tall road humps (pointed to by the cyan arrow in Fig. 5) that violate the planar motion assumption.

The content of Fig. 6 is very important as it allows us to evaluate the quality of motion estimation. In this figure, we show a textured top viewed 2D reconstruction of the whole path. Observe that this image is not an aerial image but

is an image mosaicing. Every input image of this mosaic was obtained by an Inverse Perspective Mapping (IPM) of the original omnidirectional image onto an horizontal plane. After being undistorted through IPM, these images have been merged together using the 2D poses estimated by our visual odometry algorithm. The estimated trajectory of the camera is shown superimposed with red dots. If you visually and carefully compare the mosaic (Fig. 6) with the corresponding aerial image (Fig. 5), you will recognize in the mosaic the same elements that are present in the aerial image, that is, trees, white footpaths, pedestrian crossings, roads' placement, etc. Furthermore, you can verify that the location of these elements in the mosaic fits well the location of the same elements in the aerial image.

As we mentioned already in Section 3, the fact of reducing the field of view allows us to reduce most of the problems introduced by the translation, like sudden changes in parallax. We found that the best performance in terms of closeness to the ground truth is obtained when  $FOV=10\ deg$ . This choice was a good compromise between accuracy and sensitivity to calibration errors.

## 6 Conclusion

In this paper, we described an algorithm for computing the ego-motion of a vehicle relative to the road under planar motion assumption. As only input, the algorithm uses images provided by a single omnidirectional camera. The front ends of the system are two different methods. The first one is a pure feature based method that implements the standard Triggs algorithm to compute the relative motion between two frames. The second one is an appearance based approach which gives high resolution estimates of the rotation angle of the vehicle. Using the first method to compute the vehicle displacement in the heading direction and the second one to compute the vehicle rotation proved to give very good visual odometry estimates against the pure feature based standard method.

The proposed algorithm was successfully applied to videos from an automotive platform. We gave an example of camera trajectory estimated purely from omnidirectional images over a distance of 400 meters. For performance evaluation, the estimated path was superimposed onto a aerial image of the same test environment and a textured 2D reconstruction of the path was done.

## Acknowledgment

The research leading to these results has received funding from the European Commission Division FP6-IST Future and Emerging Technologies under the contract FP6-IST-027140 (BACS). The authors would also like to say thanks to Dr. Pierre Lamon and Luciano Spinello for their useful helps, suggestions, and discussions.

## References

1. Moravec, H.: Obstacle Avoidance and Navigation in the Real World by a Seeing Robot Rover. PhD thesis, Stanford University (1980)
2. Jung, I., Lacroix, S.: Simultaneous localization and mapping with stereovision. In: *Robotics Research: The 11th International Symposium* (2005)
3. Nister, D., Naroditsky, O., Bergen, J.: Visual odometry for ground vehicle applications. *Journal of Field Robotics* (2006)
4. Maimone, M., Cheng, Y., Matthies, L.: Two years of visual odometry on the mars exploration rovers: Field reports. *Journal of Field Robotics* 24(3), 169–186 (2007)
5. Corke, P.I., Strelow, D., Singh, S.: Omnidirectional visual odometry for a planetary rover. In: *IROS* (2004)
6. Wang, H., Yuan, K., Zou, W., Zhou, Q.: Visual odometry based on locally planar ground assumption. In: *IEEE International Conference on Information Acquisition*, pp. 59–64 (2005)
7. Ke, Q., Kanade, T.: Transforming camera geometry to a virtual downward-looking camera: Robust ego-motion estimation and ground-layer detection. In: *CVPR 2003* (June 2003)
8. Guerrero, J.J., Martinez-Cantin, R., Sagues, C.: Visual map-less navigation based on homographies. *Journal of Robotic Systems* 22(10), 569–581 (2005)
9. Liang, B., Pears, N.: Visual navigation using planar homographies. In: *IEEE ICRA*, pp. 205–210 (2002)
10. Lowe, D.: Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision* 20, 91–110 (2003)
11. Tsai, R., Huang, T.: Estimating three-dimensional motion parameters of a rigid planar patch. *IEEE Trans. Acoustics, Speech and Signal Processing* 29(6), 1147–1152 (1981)
12. Longuet-Higgins, H.C.: The reconstruction of a plane surface from two perspective projections. *Royal Society London* 277, 399–410 (1986)
13. Faugeras, O.D., Lustman, F.: Motion and structure from motion in a piecewise planar environment. *International Journal of Pattern Recognition and Artificial Intelligence* (3), 485–508 (1988)
14. Fischler, M.A., Bolles, R.C.: Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM* 24(6), 381–395 (1981)
15. Scaramuzza, D., Martinelli, A., Siegwart, R.: A flexible technique for accurate omnidirectional camera calibration and structure from motion. In: *ICVS* (january 2006)
16. Labrosse, F.: The visual compass: performance and limitations of an appearance-based method. *Journal of Field Robotics* 23(10), 913–941 (2006)
17. Wunderlich, W.: Rechnerische Rekonstruktion eines ebenen Objekts aus zwei Photographien. *Mitteilungen der geodätischen Institute, TU Graz* 40, 365–377 (1982)
18. Triggs, B.: Autocalibration from planar scenes. In: Burkhardt, H.-J., Neumann, B. (eds.) *ECCV 1998. LNCS*, vol. 1406, pp. 89–105. Springer, Heidelberg (1998)