

Efficient Probabilistic Planar Robot Motion Estimation Given Pairs of Images

Olaf Booij and Ben Kröse
University of Amsterdam, The Netherlands
olaf.booij@xs4all.nl, b.j.a.krose@uva.nl

Zoran Zivkovic
NXP Research, the Netherlands
zoran.zivkovic@nxp.com

Abstract—Estimating the relative pose between two camera positions given image point correspondences is a vital task in most view based SLAM and robot navigation approaches. In order to improve the robustness to noise and false point correspondences it is common to incorporate the constraint that the robot moves over a planar surface, as is the case for most indoor and outdoor mapping applications. We propose a novel estimation method that determines the full likelihood in the space of all possible planar relative poses. The likelihood function can be learned from existing data using standard Bayesian methods and is efficiently stored in a low dimensional look up table. Estimating the likelihood of a new pose given a set of correspondences boils down to a simple look up. As a result, the proposed method allows for very efficient creation of pose constraints for vision based SLAM applications, including a proper estimate of its uncertainty. It can handle ambiguous image data, such as acquired in long corridors, naturally. The method can be trained using either artificial or real data, and is applied on both controlled simulated data and challenging images taken in real home environments. By computing the maximum likelihood estimate we can compare our approach with state of the art estimators based on a combination of RANSAC and iterative reweighted least squares and show a significant increase in both the efficiency and accuracy.

I. INTRODUCTION

Various vision based topological mapping [1, 2], view based geometrical mapping [3, 4] and robot navigation [5] approaches are based on the ability to compare pairs of images. A common way to do this is to automatically find similar looking image points [6], as done for two panoramic images in Figure 1. Because part of these point correspondences are the projections of the same 3D landmarks in the environment, they can be used to determine the relative camera pose up to an unknown scale [7]. A major challenge in determining the relative pose given point correspondences, is that a large percentage does not correspond to the same 3D landmark, but are so called *mismatches*. In addition the image point locations of correct matches are noisy, caused by for example noise of the imaging device and errors in the calibration.

To cope with this, so called *robust* algorithms are needed. There are three robust methods commonly used: RANSAC (RANdom SAMple Consensus) [8], M-Estimators (Maximum likelihood Estimators) [9] and the Hough Transform [10]. State of the art relative pose estimators first use RANSAC combined with the closed form Eight [11] or Five point algorithm [12] for an initial estimate and then apply robust iterative reweighted least squares (IRLS) techniques

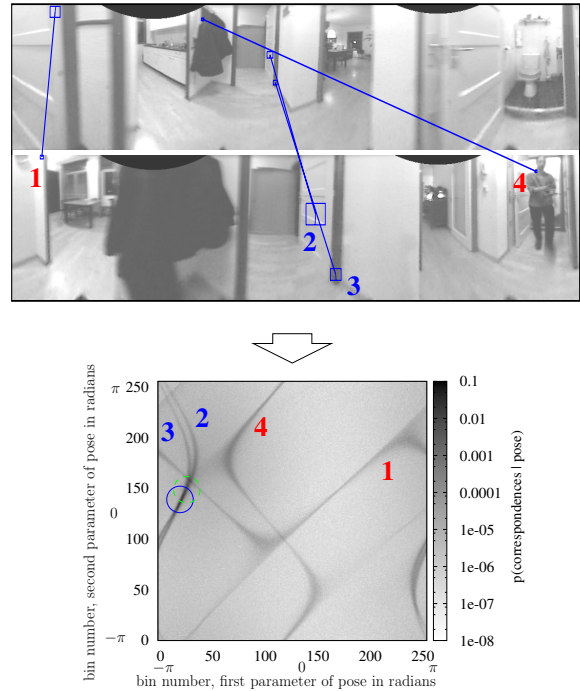


Fig. 1. An example. Applying standard SIFT matching on two panoramic images resulted in only four point correspondences including 2 mismatches (red numbers) and 2 almost degenerate correct matches (blue numbers). Still, the proposed method computes a full likelihood over the different possible relative robot poses, and the maximum likelihood (solid blue circle) is close to the ground truth (dashed green circle). The numbers relate the different correspondences in the image pair with the curves in the solution space.

such as M-Estimators in combination with the Eight point algorithm [11] to improve it. Both RANSAC and M-Estimators try to find a maximum likelihood solution by first rejecting mismatches and very noisy correspondences using an error threshold and base a least square solution on the remaining matches. Sophisticated techniques in Computer Vision try to determine this threshold from the image data itself [13, 14, 15]. In the field of Robotics, where characteristics of the camera are available, this threshold is usually determined through some calibration procedure. The Hough Transform on the other hand, if seen probabilistically, computes the full likelihood on a discrete grid of poses without making an explicit distinction between inliers and outliers. Because space requirements grow

exponential with the number of parameters, it is in general not suited for pose estimation problems, although combinations of the Hough Transform with RANSAC [16] do exist, as well as methods that treat rotation and translation estimation separately [17, 18].

An approach to make pose estimation easier is to incorporate constraints on the possible relative poses. If the robot drives over a planar surface, then the camera can only rotate around a certain fixed axis which is perpendicular to the two dimensional translation direction. Given that the scale cannot be determined, the number of degrees of freedom for this planar relative pose problem reduces to two. This constraint can be used to improve indoor mapping application using wheeled robots [5], but also vision based outdoor mapping using vehicles driving over planar roads [1, 19, 20].

Commonly this constraint is imposed rather heuristically, for example by taking only the horizontal displacement of image points into account [4, 20]. A more proper solution is proposed by Brooks [21] which formulates a least square approach to the planar relative pose given noise free correspondences. This result was used in [22, 23, 2] for various robotics applications all in combination with RANSAC. In [22, 24] it was shown that two correspondences are enough to solve the problem and both suggest algorithms, which are briefly evaluated in combination with RANSAC. Makadia proposed a direct Hough based method to estimate planar relative poses [25]. However, it also used a Spherical Fourier Transform, which made it impractical due to its large computational costs.

The key concept of RANSAC based algorithms is to evaluate the solution space for randomly sampled poses. Discretizing and analyzing the whole solution space using a Hough-like approach potentially leads to much more robust results. In this paper we show that this is feasible for the planar motion estimation case without the additional approximations as in [17, 18]. Our results demonstrate the greatly increased robustness over the random sampling techniques.

Furthermore, we present a probabilistic representation where we learn the needed conditional distributions from real or simulated data without any parameters, except the size of the discretization grid. This allows fast tuning of the approach for different robots and cameras since most standard approaches have a set of parameters that need to be tuned [4, 3]. Additionally, the learned conditional distributions properly capture other sources of uncertainty that are difficult to explicitly account for, like the shaking of the robot and inaccurate camera calibration. Therefore, more accurate and robust results are obtained as demonstrated in the experimental section. Finally, testing the whole solution space potentially can be computationally expensive. We present an efficient implementation using a precomputed look up table (LUT). The geometry of the problem is analyzed and a parameterization is proposed to reduce the dimensionality of the needed LUT to only three dimensions. The experimental results show that the whole scheme can be even faster than the common sampling based approach while at the same time we get more robust and accurate results.

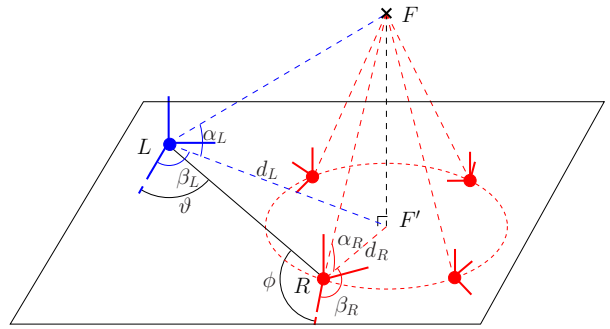


Fig. 2. 3D visualization of two cameras, L and R , positioned on the same plane both observing a landmark F . The dashed circle on the ground plane indicates the possible positions for robot R given the pose of L and observations of F .

The rest of the paper is organized as follows. First, in Section II we formalize the planar relative pose problem and describe how noise free correspondences relate to it. This relation is used in Section III to derive the novel estimator that can be trained using real image data. In Section IV we apply it on both simulated data and image data-sets taken in real home environments and compare the Maximum Likelihood estimates with the results from a planary constrained RANSAC combined with an M-Estimator. In Section V we discuss the qualitative advantage of having a full likelihood solution. Finally, in Section VI, we draw conclusions and propose some directions for improvement.

II. RELATING CORRESPONDENCES TO POSES

The planar relative pose, which can be seen as a 2D translation and rotation, minus scale, can be parameterized in different ways. We choose to parameterize it using two angles ϑ and ϕ , see Figure 2. Angle ϑ denotes the direction of the translation, or heading, to robot R in the coordinate frame of robot L and angle ϕ denotes the heading to robot L in the coordinate frame of robot R . Another common parameterizations is using the heading and the rotation of robot R in the frame of robot L , such as in [21, 22]. However, our parameterization nicely reflects the symmetry of the problem.

It is, for now, assumed that image point correspondences are obtained by a noise free projection of landmarks, without mismatches. An image point is usually denoted by a 3D vector of unit length $\mathbf{x} = [x, y, z]'$, where the z -axis is pointing in the direction of the camera axis and the y -axis is in the planar case perpendicular to the ground plane. Here we denote an image point by its horizontal angle $\beta = \text{atan2}(z, x)$ and its vertical angle $\alpha = \arcsin(y)$. This is similar to the azimuth and elevation used in the inverse depth parameterization by landmark based SLAM methods [26]. If two robots L and R observe a landmark F , the point correspondence is thus denoted by α_L, β_L and α_R, β_R .

The 3D problem as shown in Figure 2 can be reduced to a simpler 2D problem, by projecting F on the plane getting F' . We define d_L and d_R as the distances of L to F' and R to F'

respectively. The length of F to F' can now be expressed by:

$$\frac{\overline{FF'}}{d_L} = \frac{\tan(\alpha_L)}{d_L} = \frac{\tan(\alpha_R)}{d_R}, \quad (1)$$

which results in the following ratio r of d_L and d_R :

$$r = \frac{d_L}{d_R} = \frac{\tan(\alpha_R)}{\tan(\alpha_L)}. \quad (2)$$

Angle $\angle F'RL$ can then be found using the Law of Sines:

$$\frac{\sin(\angle F'RL)}{d_L} = \frac{\sin(\beta_L - \vartheta)}{d_R},$$

$$\angle F'RL = \arcsin\left(\frac{d_L}{d_R} \sin(\beta_L - \vartheta)\right). \quad (3)$$

Angle ϕ can now be expressed as a function of angle ϑ and a single point correspondence by adding the horizontal observation angle β_R , and using Equation (2):

$$\phi = \beta_R + \arcsin\left(\frac{\tan(\alpha_R)}{\tan(\alpha_L)} \sin(\beta_L - \vartheta)\right). \quad (4)$$

We could also rewrite this formula into:

$$\vartheta = \beta_L + \arcsin\left(\frac{\tan(\alpha_L)}{\tan(\alpha_R)} \sin(\beta_R - \phi)\right), \quad (5)$$

which is evident given the symmetry between ϑ and ϕ .

We now obtained functions that map ϑ to ϕ and vice versa given a single point correspondence. These can be used to plot curves of possible relative robot poses as shown in Figure 3(a) for some randomly picked point correspondences. Looking closely at this figure one can see that some pairs of curves intersect each other twice. This shows that, although it was assumed that two point correspondences can be used to solve the relative pose problem [22, 24], in some cases two different relative robot poses result in the same two point correspondences.

Based on Equation 4 an algorithm was derived that computes both these solutions given two correspondences, which can be combined with hypothesize and test schemes such as RANSAC. The exposition of this ‘‘Planar Two point algorithm’’ is rather tedious and described in [27]. For completeness it is taken into account in the experiments.

III. FULL LIKELIHOOD ESTIMATOR

We now use the relation between noise free correspondences and relative poses to develop an algorithm that can deal with noisy correspondences, including mismatches, see Figure 3(b). The problem can be formulated as determining the negative log likelihood of each pose (ϑ, ϕ) given n point correspondences $\{\xi_1, \dots, \xi_n\}$, each parameterized by the angles $\xi_i = (\alpha_{Li}, \beta_{Li}, \alpha_{Ri}, \beta_{Ri})$ [28]:

$$\mathcal{L}(\vartheta, \phi) = \sum_i \mathcal{L}_i(\vartheta, \phi) \quad (6)$$

$$= \sum_i -\log p(\xi_i | \vartheta, \phi), \quad (7)$$

where we use the common assumption that the correspondences are independent given the relative pose. We discretize

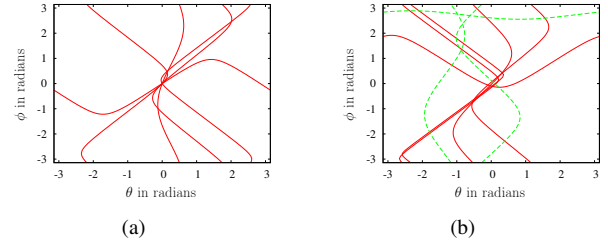


Fig. 3. Visualization of the possible relative robot poses by plotting ϕ as a function of ϑ for multiple correspondences randomly picked using a relative pose with $\vartheta = \phi = 0$. (a) Generated using 5 noise free correspondences (red/solid curves). (b) Generated using 5 noisy correspondences and 3 mismatches (green/dashed curves). Note that, by coincidence, one of the mismatches corresponds to a curve close to the actual pose at $\vartheta = \phi = 0$.

the 2D likelihood into a 2D histogram. For each bin we need to sum the contribution of all point matches. Therefore for each point match we need to calculate the corresponding 2D histogram approximating its log likelihood contribution $\mathcal{L}_i(\vartheta, \phi)$, which is computational costly. In the next section, we show how this can be efficiently performed using a precomputed look-up table.

A. Look up table

The negative log likelihood of a single correspondence ξ is given by $-\log p(\alpha_L, \beta_L, \alpha_R, \beta_R | \vartheta, \phi)$. Thus, if we would naively construct a look up table, it would have 6 dimensions, 4 for the point correspondence angles and 2 for the relative pose. In order to keep the size of the LUT comprehensible, the 6D space should be discretized in large bins, resulting in a large discretization error. Fortunately, the dimensionality of the LUT can be reduced to only 3 dimensions by using the one point mapping function introduced in Section II and some common assumptions about the noise characteristics of image points.

The function given in Equation (4) can be written as

$$\phi - \beta_R + \arcsin\left(\frac{\tan(\alpha_R)}{\tan(\alpha_L)} \sin(\vartheta - \beta_L)\right) = 0. \quad (8)$$

As can be seen the terms α_L and α_R are only used in the combination $\frac{\tan(\alpha_R)}{\tan(\alpha_L)}$. Also, variables ϑ and β_L , and variables ϕ and β_R are only used in the combinations $\vartheta - \beta_L$ and $\phi - \beta_R$ respectively, which describe the horizontal angles to the landmark relative to the heading of the cameras. The joint probability of observing a correspondence under a certain relative pose can thus be represented as:

$$p(\alpha_L, \beta_L, \alpha_R, \beta_R, \vartheta, \phi) = p\left(\frac{\tan(\alpha_R)}{\tan(\alpha_L)}, \vartheta - \beta_L, \phi - \beta_R\right). \quad (9)$$

This holds if we assume that the noise of the horizontal and vertical view angles do not depend on their value, which is similar to the common assumption that the noise of the pixel locations is homogeneous.

Because of the symmetry of the representation of the relative planar pose, the two points of the point correspondences and

the heading angles can be swapped giving the same result:

$$p(r, \vartheta - \beta_L, \phi - \beta_R) = p\left(\frac{1}{r}, \phi - \beta_R, \vartheta - \beta_L\right), \quad (10)$$

with $r = \frac{\tan(\alpha_R)}{\tan(\alpha_L)}$ for convenience. In practice this means that we only have to construct a LUT for $0 < r < 1$ and swap $\vartheta - \beta_L$ with $\phi - \beta_R$ and use $\frac{1}{r}$ instead of r if $r > 1$.

The likelihood can be determined from the joint probability by dividing by the probability of the pose $p(\vartheta, \phi)$:

$$p(\xi|\vartheta, \phi) = p(\alpha_L, \beta_L, \alpha_R, \beta_R|\vartheta, \phi) \quad (11)$$

$$= p(\alpha_L, \beta_L, \alpha_R, \beta_R, \vartheta, \phi) / p(\vartheta, \phi) \quad (12)$$

$$= p\left(\frac{\tan(\alpha_{Ri})}{\tan(\alpha_{Li})}, \vartheta - \beta_{Li}, \phi - \beta_{Ri}\right) / p(\vartheta, \phi) \quad (13)$$

During the construction the LUT, one can take care that the different relative poses are uniformly distributed as described below. This makes the likelihood proportional to the joint:

$$p(\xi|\vartheta, \phi) \propto p\left(\frac{\tan(\alpha_R)}{\tan(\alpha_L)}, \vartheta - \beta_L, \phi - \beta_R\right). \quad (14)$$

Efficiently determining a full likelihood over the relative poses given a set of correspondences is now straightforward. For each correspondence ξ_i we compute the value of $\frac{\tan(\alpha_{Ri})}{\tan(\alpha_{Li})}$ and pick the corresponding 2D slice of the look up table. Then we shift it in the direction of β_{Li} and β_{Ri} , wrapping the values at the borders. This results in the negative log likelihood of each pose given a correspondence, which can be summed for the different correspondences, resulting in a full likelihood. Algorithm 1 summarizes this procedure.

Algorithm 1 Full likelihood estimator using a LUT

HIST \leftarrow 2d array of zeros

for all correspondences $\{\alpha_L, \beta_L, \alpha_R, \beta_R\}$ **do**

 compute $r = \frac{\tan(\alpha_R)}{\tan(\alpha_L)}$

 HIST $+=$ LUT($d(r)$) shifted by $d(\beta_L)$ and $d(\beta_R)$

 where $d(\cdot)$ denotes discretization

B. Constructing a look up table

The LUT representing the negative log likelihood can be constructed from existing data. This data can be generated by a simulator modeling the planar pose problem including a vision system. Better is to use a representative image set for which ground truth robot pose data is available. The main problem of real data is that the relative poses are in general not uniformly distributed. This invalidates the simplification proposed in Equation (14) and results in a bias of the LUT towards certain poses which are overrepresented in the dataset. However, we can easily compensate for this problem.

When constructing the LUT, we explicitly take the probability of the relative pose into account (see Equation (13)). In a first step, a 2D discretized probability $p(\vartheta, \phi)$ is constructed by making a histogram for all poses in the dataset and normalizing it. Then, in a second step, the dataset is used to build the 3D LUT like for the simulator, with the difference that for

each pose correspondence it adds $\frac{1}{p(\vartheta, \phi)}$ to the 3D histogram. Again, each value of the histogram is replaced by its negative log, resulting in a proper LUT. Algorithm 2 summarizes the procedure to build a LUT. To keep the algorithm comprehensible, it ignores the symmetry between ϑ and ϕ .

A second problem, which can not be circumvented, is that the amount of data in an image dataset is limited. As a consequence we usually can not construct a LUT with a very high number of bins. In the next section we evaluate, among other things, the consequences of such a smaller LUT.

Algorithm 2 Constructing a LUT from real data

PP \leftarrow 2d array of zeros

for all train data $\{\vartheta, \phi, \alpha_L, \beta_L, \alpha_R, \beta_R\}$ **do**

 PP($d(\vartheta)$, $d(\phi)$) $+=$ 1

LUT \leftarrow 3d array of zeros

for all train data $\{\vartheta, \phi, \alpha_L, \beta_L, \alpha_R, \beta_R\}$ **do**

 compute $r = \frac{\tan(\alpha_R)}{\tan(\alpha_L)}$

 LUT($d(r)$, $d(\vartheta - \beta_L)$, $d(\phi - \beta_R)$) $+=$ 1/PP($d(\vartheta)$, $d(\phi)$)

LUT \leftarrow - log each entry of LUT

IV. EXPERIMENTS

By determining the Maximum Likelihood from the estimated full likelihood (LUT ML), we can compare our method for planar relative pose estimation with the state of the art robust methods using RANSAC combined with an M-Estimator as described briefly in Section I. We combined these with the Planar Two point, the Planar Three point and the general Eight point algorithm. We used simulated data and 4 datasets obtained from a robot with an omnidirectional vision system. The methods are compared on the basis of their robustness against mismatches and noise.

We evaluate the estimated heading and rotation angle as in [29] by taking the absolute difference with the ground truth values. Because these errors are not normally distributed we use the median to describe the error distribution. A robust way to describe the spread of these medians is to use the Median of Absolute Deviations (MAD). Another important evaluation criterion is the computational time used by the algorithms. For all experiments we report these times, as implemented in C++ and run on the single 2 Ghz CPU core of a Pentium PC.

A. Experiments on simulated data

Using simulated data allows us to control the projection noise and number of mismatches.

1) *Data*: Data was simulated by randomly picking a uniformly distributed point cloud of 3D landmarks inside a sphere of size 2 around the origin. Two random camera poses on a circle with radius 1 in the x-y plane around the origin are chosen. From these the ground truth values for ϑ and ϕ are determined. Note that the distribution of ϑ and ϕ is approximately uniform.

A set of point correspondences is constructed by projecting the landmarks on a spherical shaped image surface with a

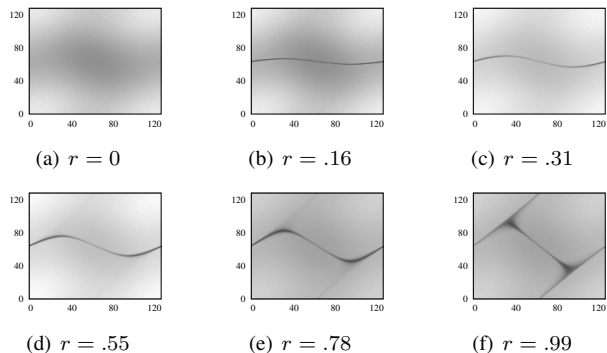


Fig. 4. Look up table obtained from simulated data as described in the experiment section.

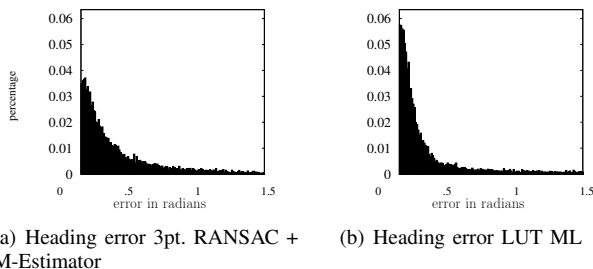


Fig. 5. Comparison of the RANSAC+M-Estimator combined with the Three point algorithm and the proposed LUT ML method for 90% mismatches. The distribution of rotation errors shows a similar pattern.

radius of 1 around the camera pose. Thus, an ideal omnidirectional camera model is used with a full 360 degrees view angle in the horizontal and vertical direction. An amount of normally distributed noise with zero mean and standard deviation 0.01 is added to these projections. In addition, mismatches are added by creating false correspondences between projections of different landmarks. We use a mismatch rate of 90%.

2) *Setup*: A LUT was constructed using the same simulator as described above. The number of bins to represent ϕ , ϑ and r were all 128, which caused the computational time to be comparable with that of the RANSAC+M-Estimator. The number of point correspondences used was 10^{10} , which took 3 hours to build. The error threshold for both RANSAC and the M-Estimator were set according to the projection noise of the simulator. The number of RANSAC iterations was set to 100.

3) *Resulting look up table*: Figure 4 shows the resulting look up table. In 4(a) one can see that point correspondences with a r value close to zero, do not tell that much about the data. This is due to the fact that there is a high chance that it resulted from a mismatch. Note also that the histogram corresponding to a r value close to 1 in 4(f) is almost symmetric. Indeed if r is exactly 1 then $r = \frac{1}{r}$ and we could swap ϕ and θ .

4) *Resulting distributions*: The distribution of errors is given in Figure 5. The LUT ML method results is superior for 90% percent mismatches. Also, both distributions have long

tails, indicating the need of robust statistics.

5) *Sensitivity to mismatches*: To test the robustness of the different methods to mismatches, we vary the number of mismatches, from 50% to 99%. In total 10^5 iterations were conducted. In Figures 6(a) and 6(e) the resulting estimation errors are plotted.

Both the heading and rotation results show the same trend. The error of the RANSAC+M-Estimator with the Eight point algorithm, which does not take planarity into account, increases fast if more than half of the correspondences are mismatches. The RANSAC+M-Estimator with the Two point algorithm is better, but is in turn outperformed by the Three point algorithm which starts diverging at 70% mismatches. The accuracy of the LUT ML estimator starts diverging at 85% and seems particularly robust against high mismatch rates. For mismatch rates under 60% the relative error is somewhat higher which can be explained by the discretization process used in the algorithm.

6) *Sensitivity to violations of the planar assumption*: In practice the motion of a robot is never strictly planar. Therefore, we test the behavior of the algorithms when small rotations around the x-axis, and around the z-axis, are added to the poses. We learn two separate LUTs for this experiment: one as described above using 10^7 simulated correspondences named LUT con, and one in which we additionally simulate small rotations around the x-axis, and around the z-axis, both up to .25 radians, named LUT rot. In Figures 6(b) and 6(f) the median errors are plot against the amount of pitch and roll, for a mismatch rate of 60%. As can be seen at an angle of .3 radians the RANSAC method and conventional LUT ML method have the same or worse error as the RANSAC with Eight point algorithm which is not influenced by non planar motion. Note that this corresponds to a percentage-grade of 31%. The LUT ML method with LUT con, learned to be more robust and is better than the RANSAC with Eight point algorithm up to .4 radians.

B. Experiments on real data

We compared the performance of our method with other methods on more than $3 * 10^6$ real image pairs.

1) *Data*: We used three distinct image datasets. The first two were obtained using our Nomad Super Scout II and the third by the ‘Biron’ robot from the University of Bielefeld. On both an omnidirectional vision system was mounted consisting of a conventional Firewire camera pointing upwards to a hyperbolic mirror. Note that the proposed methods do not require omnidirectional images. All sets are taken in real home environments. In the ‘Almere 4’ set there are some people walking in the rooms, the ‘Spaan 1’ set is taken during evening hours, and the ‘Biron 1’ set is taken in a feature poor home. The ground truth robot poses for the home sets were obtained by applying the SLAM algorithm described in [30] on laser scans and odometry. The home sets are available for research purposes, see [31] for details.

2) *Setup*: From every dataset we use every pair of images. We discard the images taken at the same position. Also, if

images are taken at more than 5 meters apart for the Almere 4 set or more than 3 meters apart for the other sets, then the chance of finding point correspondences is small, so we also discard these pairs. Still, for each set there are around 10^6 image pairs left.

To extract point correspondences from the image pairs, the SIFT algorithm is used [6]. First omnidirectional images are mapped to panoramic images [32], from which the SIFT feature points are found. On average 238 features were found per image. These features are described by the standard SIFT descriptor of 128 dimensions. If two features in the same image have a small distance in descriptor space then they are removed. A set of point correspondences between two images is determined by applying the standard matching scheme as described in [6]. This resulted in on average 25 matches per image pair. The groundtruth relative pose was computed from the groundtruth robot positions.

3) Sensitivity to mismatches, trained with simulated data:

We first use a LUT constructed using the simulator. In order to evaluate the performance of the methods we would like to vary the number of mismatches. This can not be controlled in real data, therefore we made subsets of the data on the basis of the distance between the poses. We assume that for larger distances it is more difficult to find matching features. In Figure 6(c) and 6(g) the heading and rotation error of the different methods is plot as a function of the distance between the images for dataset ‘Almere 4’. It is clear that on a whole the errors are much larger than was the case for the simulation data. This is partly due to the fact that some of the views were obstructed by furniture, walls or people walking in the environment.

In the plot of the heading error (Figure 6(c)) one can see that the RANSAC+M-Estimator combined with the Two point algorithm is outperformed by the three point algorithm version, which in turn is clearly outperformed by the novel LUT method for distances larger than 1.5 meters. The accuracy of the RANSAC+M-Estimator combined with the Eight point algorithm is not much worse. This could have been caused by the fact that the robot was leaning over when accelerating, slightly violating the planarity constraint. For the rotation error (Figure 6(g)), the improvement of the ML estimator over the RANSAC+M-Estimator with Three Point is less clear.

4) *Sensitivity to mismatches, trained with real data:* Next we constructed a look up table using all the image pairs of the Almere 4 set that were within a 5 meter distance. We used two different binsizes, the first had 128 bins for all three dimensions and the other 16. Although about 10^6 image pairs were used, the LUT with 128^3 bins was not as smooth as the one learned with simulated data. The proposed method was applied on the Spaan 1 set. The Maximum Likelihood solutions were compared to the RANSAC+M-Estimator combined with the Three point algorithm and solutions given two LUTs based on the simulator, also with 128^3 and 16^3 .

In Figure 6(d) and 6(h) the results are shown. As can be seen the overall accuracy is less than for the Almere 4. A reason for this could be the motion blur, caused by the bad

TABLE I
AVERAGE COMPUTATIONAL TIME USAGE PER RELATIVE POSE ESTIMATE
IN MILLISECONDS FOR THE DIFFERENT METHODS.

ML				RANS+M-Est		
128	64	32	16	8pt	3pt	2pt
1.3	0.28	0.07	0.036	3.6	3.8	0.68

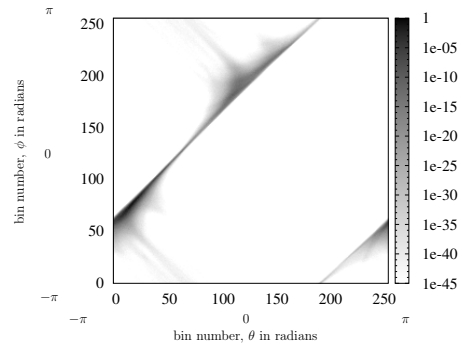


Fig. 8. The log posterior computed from the point correspondences between image 857 and 897 of Almere set 4. The distance between these camera positions was about 1 cm, while the rotation was 90 degrees.

illumination of this dataset. The LUT ML with 128^3 bins based on the simulator performs best, followed by the 128^3 bins LUT based on the real images. Probably this somewhat worse performance is due to the limited number of training data. For the much smaller LUTs with 16^3 bins this seems to be less problematic, visible by the improvement of the LUT based on real images over the one based on simulated data.

5) *Averages over the data sets:* Application on other datasets and different bin sizes resulted in comparable errors. Figure 7 summarizes these results. Note that all three RANSAC based methods failed to robustly estimated poses for the difficult Biron 1 set as opposed to the proposed method.

6) *Binsize vs CPU time:* To evaluate the influence of different binsizes for the look up table, we tested the ML method for different numbers of bins. In Table I the average computational time in milliseconds is given per image pair for the ‘Almere 4’ set. Other datasets showed similar trends for both the LUTs constructed using the simulator and the real data. As can be seen small look up tables result in a large speed up, but this comes at the cost of more error (see Table 7). The RANSAC+M-Estimator combined with the Planar Three point algorithm is three times slower than the LUT method with 128^3 bins.

V. DISCUSSION

An important advantage of the LUT based pose estimator is that it provides a full likelihood over the discretized space of possible relative poses. Thus, apart from finding a Maximum Likelihood solution, as shown in the Experiments section, this could make the method useful for a range of other applications.

A nice illustration of the usefulness of a full likelihood can be seen in Figure 8. It shows an example likelihood for a

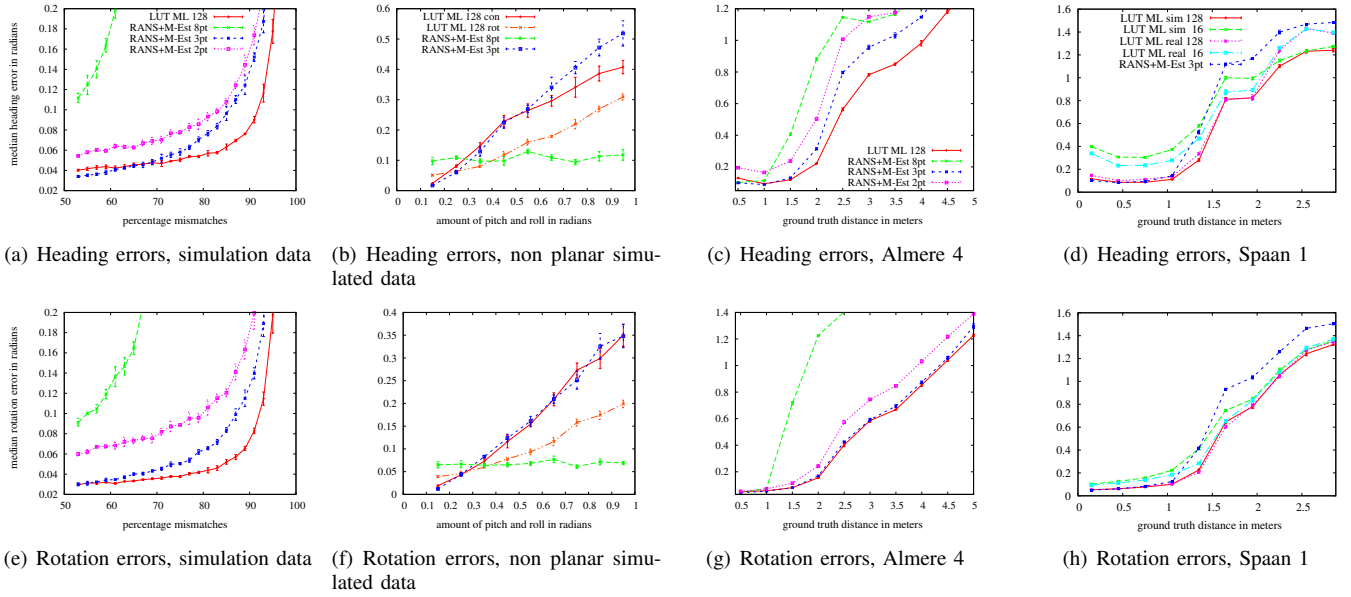


Fig. 6. Comparison of LUT ML, trained using a simulator or real images and RANSAC+M-Estimator combined with different algorithms on the simulation dataset for different number of mismatches and amount of planar motion violation and Almere 4 and Spaan 1 for different distances between the image pairs. The MAD is used to draw confidence intervals of the medians.

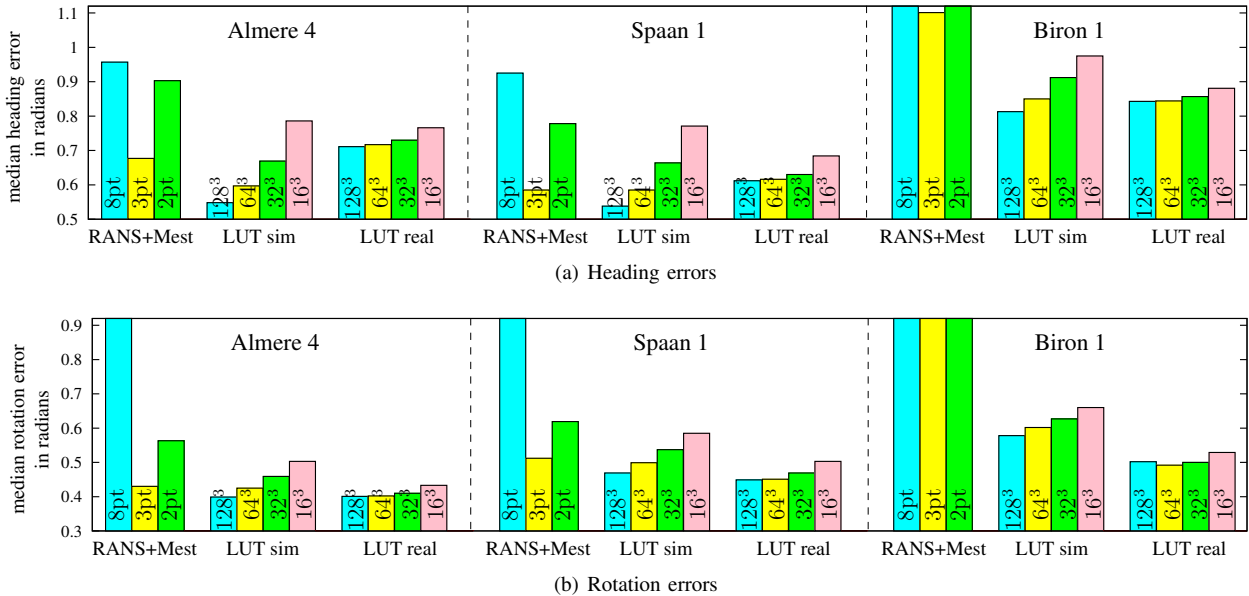


Fig. 7. (a) Heading errors and (b) rotation errors for the different estimators applied on all three datasets.

typical situation that occurred in the Almere 4 set. In this case the robot did not move forward but rotated on the spot. Thus, the heading of the robot, ϑ , can not be determined. This is correctly reflected by the estimated posterior. The rotation, on the other hand, can be determined. This can be seen by the diagonal relationship between ϑ and ϕ in the posterior.

The proposed method could be readily applied on particle filter based robot localization schemes [33] where each hypothesized robot pose can be weighed by the likelihood given newly acquired images. Also, geometric View based

SLAM [3, 4] could benefit from the proposed method, because the uncertainty of the Maximum Likelihood can be estimated easily from the full likelihood. For example by fitting a Von Mises or mixture of Von Mises distributions on the discretized likelihood space [28].

Another task that is very much suited for the proposed example is that of topological mapping. Some state of the art topological mapping approaches use proper probabilistic data association techniques to compare pairs of images [1]. However, in addition they apply ad hoc rules to check whether the

matched point correspondences fit in a certain local geometry, estimated using RANSAC. Because of the probabilistic nature of the proposed method, it is straightforward to combine it with these proper data association techniques, ending up in a fully probabilistic topological mapping method.

VI. CONCLUSION

In this paper we propose a novel approach to solve planar relative pose estimation from image point correspondences. We have shown the advantage of discretizing and analyzing the complete solution space, which is in the planar motion case 2 dimensional. Probabilistic methods were proposed that learn the likelihood over this space from a training set of representative images. Experiments on challenging image sets acquired in real homes showed a 20% increase in accuracy with respect to state of the art methods consisting of a planar constrained RANSAC and M-Estimators.

In addition an efficient technique was presented for building a concise look up table of the likelihood, reducing the estimation process to simple look ups. Computing a full likelihood given two images costs as little as 36 microsecond, as compared to the 3 milliseconds RANSAC uses. This could even be improved upon, for example, by using a multi-resolution approach as described in [34], in which a small look up table is used to isolate candidate areas for the ML solution, which are then investigated further using a bigger look up table. Another possibility is implementing the method on a GPU, which can much more quickly manipulate 2D histograms.

Continuing research will focus on taking advantage of the full likelihood that is estimated by the method. We foresee improvements in both geometric SLAM as in topological mapping, especially when it comes to uncertainty estimation.

ACKNOWLEDGMENTS

We would like to thank Gwenn Englebienne and Leo Dorst for fruitful discussions and contributions to this paper.

REFERENCES

- [1] M. Cummins and P. Newman, "Highly scalable appearance-only SLAM - FAB-MAP 2.0," in *Proceedings of Robotics: Science and Systems (RSS)*, Seattle, USA, 2009.
- [2] C. Valgren and A. J. Lilienthal, "Incremental spectral clustering and seasons: Appearance-based localization in outdoor environments," in *ICRA*. Pasadena, California, USA: IEEE, May 2008, pp. 1856–1861.
- [3] R. Eustice, "Large-area visually augmented navigation for autonomous underwater vehicles," Ph.D. dissertation, MIT - Woods Hole Oceanographic Institute, June 2005.
- [4] H. Andreasson, T. Duckett, and A. J. Lilienthal, "A minimalistic approach to appearance based visual slam," *Transactions on Robotics*, vol. 24, no. 6, pp. 991–1001, October 2008.
- [5] F. Fraundorfer, C. Engels, and D. Nister, "Topological mapping, localization and navigation using image collections," in *IROS*. San Diego, USA: IEEE/RSJ, November 2007, pp. 3872–3877.
- [6] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *Int. Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [7] R. Hartley and A. Zisserman, *Multiple view geometry in computer vision, second edition*. Cambridge University Press, 2003.
- [8] M. A. Fischler and R. C. Bolles, "Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography," *Com. of the ACM*, vol. 24, no. 6, 1981.
- [9] P. H. S. Torr and D. W. Murray, "The development and comparison of robust methods for estimating the fundamental matrix," *Int. Journal of Computer Vision*, vol. 24, no. 3, pp. 271–300, 1997.
- [10] P. V. C. Hough, "Method and means for recognizing complex patterns," 1962, u.S. Patent 3,069,654.
- [11] H. C. Longuet-Higgins, "A computer algorithm for reconstructing a scene from two projections," *Nature*, vol. 293, pp. 133–135, Sept. 1981.
- [12] D. Nistér, "An efficient solution to the five-point relative pose problem," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 26, no. 6, pp. 756–777, 2004.
- [13] B. C. Matei, "Heteroscedastic errors-in-variables models in computer vision," Ph.D. dissertation, Rutgers University, 2001.
- [14] H. Wang and D. Suter, "Robust adaptive-scale parametric model estimation for computer vision," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 26, no. 11, pp. 1459–1474, 2004.
- [15] P. H. S. Torr and A. Zisserman, "MLESAC: a new robust estimator with application to estimating image geometry," *Comput. Vis. Image Underst.*, vol. 78, no. 1, pp. 138–156, 2000.
- [16] R. den Hollander and A. Hanjalic, "A combined RANSAC-Hough transform algorithm for fundamental matrix estimation," in *18th British Machine Vision Conference*. University of Warwick, UK, 2007.
- [17] D. J. Heeger and A. D. Jepson, "Subspace methods for recovering rigid motion i: algorithm and implementation," *Int. J. Comput. Vision*, vol. 7, no. 2, pp. 95–117, 1992.
- [18] A. Censi and S. Carpin, "HSM3D: Feature-less global 6DOF scan-matching in the Hough/Radon domain," in *ICRA*. Kobe, Japan: IEEE, May 2009.
- [19] M. Pollefeys, D. Nistér, J. M. Frahm, A. Akbarzadeh, P. Mordohai, B. Clipp, C. Engels, D. Gallup, S. J. Kim, P. Merrell, C. Salmi, S. Sinha, B. Talton, L. Wang, Q. Yang, H. Stewénius, R. Yang, G. Welch, and H. Towles, "Detailed real-time urban 3d reconstruction from video," *Int. J. Comput. Vision*, vol. 78, no. 2-3, pp. 143–167, 2008.
- [20] M. Milford and G. Wyeth, "Mapping a suburb with a single camera using a biologically inspired SLAM system," *IEEE Transactions on Robotics*, vol. 24, no. 5, pp. 1038–1053, 2008.
- [21] M. Brooks, L. de Agapito, D. Huynh, and L. Baumela, "Towards robust metric reconstruction via a dynamic uncalibrated stereo head," *Image Vision Comput.*, vol. 16, no. 14, pp. 989–1002, 1998.
- [22] D. Ortín and J. M. M. Montiel, "Indoor robot motion based on monocular images," *Robotica*, vol. 19, no. 3, pp. 331–342, 2001.
- [23] J. Kosecká, F. Li, and X. Yang, "Global localization and relative positioning based on scale-invariant keypoints," *Robotics and Autonomous Systems*, vol. 52, no. 1, pp. 27–38, 2005.
- [24] T. Goedemé, T. Tuytelaars, G. Vanacker, M. Nuttin, and L. V. Gool, "Feature based omnidirectional sparse visual path following," in *IROS*. Edmonton, Canada: IEEE/RSJ, August 2005, pp. 1003–1008.
- [25] A. Makadia, D. Gupta, and K. Daniilidis, "Planar ego-motion without correspondences," in *WACV/Motion*. Los Alamitos, CA, USA: IEEE Computer Society, January 2005, pp. 160–165.
- [26] J. Civera, A. J. Davison, and J. M. M. Montiel, "Inverse depth parametrization for monocular SLAM," *IEEE Transactions on Robotics*, vol. 24, no. 5, pp. 932–945, October 2008.
- [27] O. Booij and Z. Zivkovic, "The planar two point algorithm," University of Amsterdam, Tech. Rep. IAS-UVA-09-05, 2009.
- [28] C. M. Bishop, *Pattern Recognition and Machine Learning*, ser. Information Science and Statistics. Springer, 2006.
- [29] D. H. Stewénius, C. Engels, and D. D. Nistér, "Recent developments on direct relative orientation," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 60, no. 4, pp. 284–294, 2006.
- [30] J. Folkesson, P. Jensfelt, and H. Christensen, "Vision SLAM in the measurement subspace," in *ICRA*. Barcelona, Spain: IEEE, April 2005, pp. 30–35.
- [31] Z. Zivkovic, O. Booij, B. Kröse, E. Topp, and H. I. Christensen, "From sensors to human spatial concepts: an annotated dataset," *IEEE Transactions on Robotics*, vol. 24, no. 2, pp. 501–505, April 2008.
- [32] R. Bunschoten, "Mapping and localization from a panoramic vision sensor," Ph.D. dissertation, University of Amsterdam, November 2003.
- [33] H.-M. Gross and A. Koenig, "Robust omniview-based probabilistic self-localization for mobile robots in large maze-like environments," in *ICPR (3)*, 2004, pp. 266–269.
- [34] E. B. Olson, "Real-time correlative scan matching," in *ICRA*. Kobe, Japan: IEEE, May 2009.