
SIFT Based Graphical SLAM on a Packbot

John Folkesson¹ and Henrik Christensen²

¹ Massachusetts Institute of Technology, Email: johnfolk@mit.edu

² Georgia Institute of Technology, Email hic@cc.gatech.edu

Summary. We present an implementation of Simultaneous Localization and Mapping (*SLAM*) that uses infrared (*IR*) camera images collected at 10 Hz from a Packbot robot. The Packbot has a number of challenging characteristics with regard to vision based SLAM. The robot travels on tracks which causes the odometry to be poor especially while turning. The IMU is of relatively low quality as well making the drift in the motion prediction greater than on conventional robots. In addition, the very low placement of the camera and its fixed orientation looking forward is not ideal for estimating motion from the images. Several novel ideas are tested here. Harris corners are extracted from every 5th frame and used as image features for our SLAM. Scale Invariant Feature Transform, *SIFT*, descriptors are formed from each of these. These are used to match image features over these 5 frame intervals. Lucas-Kanade tracking is done to find corresponding pixels in the frames between the SIFT frames. This allows a substantial computational savings over doing SIFT matching every frame. The epipolar constraints between all matches that are implied by the dead-reckoning are used to further test the matches and eliminate poor features. Finally, the features are initialized on the map at once using an inverse depth parameterization which eliminates the delay in initialization.

1 Introduction

³ The goal of much research over the last decade has been to give robots a sense of location and direction comparable to that of humans. Humans have very limited sense of their ego motion with eyes closed. With eyes opened one can move with more confidence. Judgements of distance and orientation are greatly improved. Comparable automatic robot vision capabilities have not yet been demonstrated. Progress has been made [1], [2], [3], [4], and [5].

The problem has been given the name of vision based simultaneous localization and mapping SLAM. An issue in making maps from features seen in images is the lack of depth information in a single image [6]. Most of the

³ This work was supported by the Swedish defense agency FMV. This support is gratefully acknowledged.

SLAM methodologies estimate a maximum likelihood solution given the motion and feature measurements. This estimate typically assumes that the map and robot pose can be described by Gaussian distributions. Features seen in a single camera image correspond to real 3D objects. The probability distribution of such an object's location given the image has the shape of an infinite cone in Cartesian 3 space. This causes considerable difficulty for SLAM methods that use xyz as the feature representation. This problem has been cleverly solved [7] by utilizing a representation in terms of the bearing from the pose of the camera at the first observation and the inverse depth in that frame. The 6 parameters of the camera pose for this frame also become part of the parameterization. The result is a 9 parameter representation of the 3D point whose distribution can be well approximated by a Gaussian.

Another issue in vision based SLAM is the extraction of pixel level features from images [8], [9] and matching these with images of the same objects from different vantage points. Scale Invariant Feature Transform, SIFT, features have been used for SLAM [10], [11], [12], [13]. SIFT features have relatively weak descriptors associated with them. These descriptors have the advantage of invariance with respect to scale. They are also fairly similar for changes in viewing angle up to about 15-20 degrees. The descriptors alone can not be used to match features however and must be combined with other evidence such as epipolar constraints or the relative locations of other matching points.

2 Graphical SLAM method

The Graphical SLAM method has been presented in [14] and [15]. We refer to those papers for the mathematical details. The measurements over some time interval from motion sensors give us a probability distribution over the robot poses at the beginning and end of the interval. This is the probability of the measurements given the state. The negative log of this is the energy contribution of the dead reckoning over the interval. Thus the adjustments to the state from a number of measurements can be made based on a minimization principle on a graph of robot pose nodes.

The camera images give us additional information on the motion and additional states corresponding to features seen in the images. These additional feature states become feature nodes in our graph. Their measurements create energy terms relating them to the robot poses. So that the graph consists of state nodes for the 6 dof robot poses at the times of the camera frames and the feature nodes of the 3D points that gave rise to image features in those frames. The connections between these state nodes is through energy nodes that correspond to measurements. The main advantage of working with this graph construction is that information can be gathered more flexibly and the linearization errors are greatly reduced. Having the pose states allows one to add and remove measurements from the solution. One can change data associations after testing how the energy would be changed. As a general principle

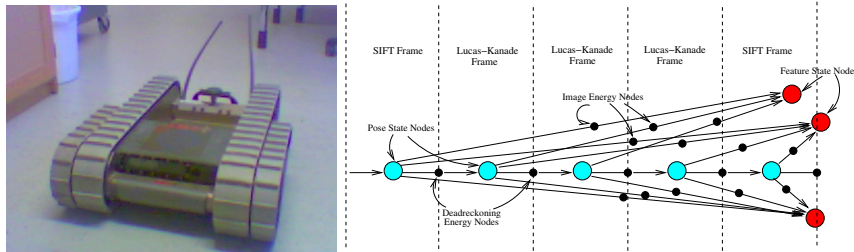


Fig. 1. The packbot robot is shown beside an illustration of a short graph segment.

one can say that it is easier to find representations that fit well to Gaussian distributions by increasing the dimension of the state space. This is true for the robot pose and, as we will see, for the feature part of the state as well.

The down side of increasing the dimensions is higher computation cost. In fact the full graphical solution gets unmanageable very quickly. However two escapes can be found. First the growth of the graph by adding more measurements does not change the solution for states far removed from the added measurements. That is unless the measurements include some new global constraints such as closing of a loop. So as long as the new measurements are not matched to older features one can simply relax the end of the graph, a constant time update. This leaves out the global localization problem completely. One is concerned with local accuracy only and the global solution is not important.

The other escape is to marginalize out part of the state from mature sections of the graph. This will allow global matching with reasonable complexity while maintaining local accuracy of the high dimensional representation. This is done by forming what we call star nodes⁴. A star node is an energy node that replaces a set of energy nodes with a linear approximation of their energies. One can then marginalize out the pose nodes that only connect to a single star node. By doing this in a local frame one obtains invariance with respect to large transformations of sections of the map.

One thus maintains the full state for the part of the graph connected to currently observed features. Then for the older section one forms a star node that marginalizes out a section of the robot path. The star node is allowed to grow until it has eliminated a predetermined number of pose nodes or has a maximum number of features. Then a new star node is begun.

3 The Feature Tracking

The SLAM algorithm is only one part of an implementation. The results will also depend on the tracking/matching of features in the sequence of im-

⁴ Exactly sparse EIF methods are similarly leave in some poses in order to facilitate the computations. These do not reduce linearization errors which star nodes do.

age frames. Harris corners have an advantage over other image features in that they have well defined pixel locations in the image. Those locations may or may not correspond to well defined 3D locations in the real world. The

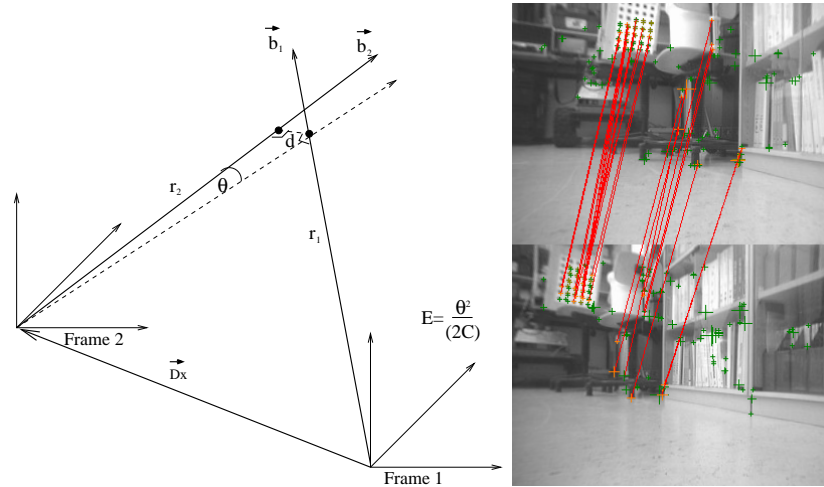


Fig. 2. On the left is the geometry of the constraint energy. Two bearings, \mathbf{b}_1 and \mathbf{b}_2 and the transformation between the camera frames, \mathbf{Dx} , is used to find θ . The energy of the constraint is $\theta^2/(2C)$, where C is the variance in θ . The two ranges r_1 and r_2 are calculated to minimize d . In the case of poorly defined ranges as when the bearings are in the \mathbf{Dx} direction we use a default distance of 5 m for r_1 . The right is one example of poor SIFT matching between frames. Here the robot motion between the images leads to significant change in the images. The nature of the scene is leading to similar SIFT descriptors for different parts of the image.

matching of the SIFT descriptors normally gives a good indication of pixel level correspondences between frames. The difficulty in this implementation was that it took about 100 msec per image to extract the Harris corners and calculate the SIFT descriptors for them. We would like to use a 10 Hz frame rate as the images can change rapidly as the robot moves about. We can only hope to achieve cpu utilizations of around 50-70% in the real time application. Thus, SIFT descriptors can not be extracted from every frame.

The solution we chose was to do Harris/SIFT extraction on every 5th frame. Then for the intervening frames we interpolate the SIFT matches on each side using Lucas-Kanade tracking [16]. This tracking maximizes the correlation to a templet of the pixels around our feature to achieve sub pixel tracking accuracy. In this way we could reduce the utilization giving us a margin for other computations.

Having tracked an image feature over 6 frames we apply the epipolar constraint to the motion in the image sequence. This we do in by utilizing our

graph of pose nodes. We form temporary constraint energy nodes connecting two poses. The energy of these constraint nodes is defined in figure 2. This is then used to find the constraint. One advantage of this definition of the constraint energy is that the epipolar constraint can be applied for features seen in the direction of the motion. Such features have no epipolar line but the constraint energy still makes sense.

These constraints are one dimensional and three of them can specify the rotation between two frames given the translation. The scale of the the translation can not be found but the energy is a function of the translation direction as well. We can differentiate the constraint energy with respect the pose states. This will give us a gradient and hessian to relax the graph with. We

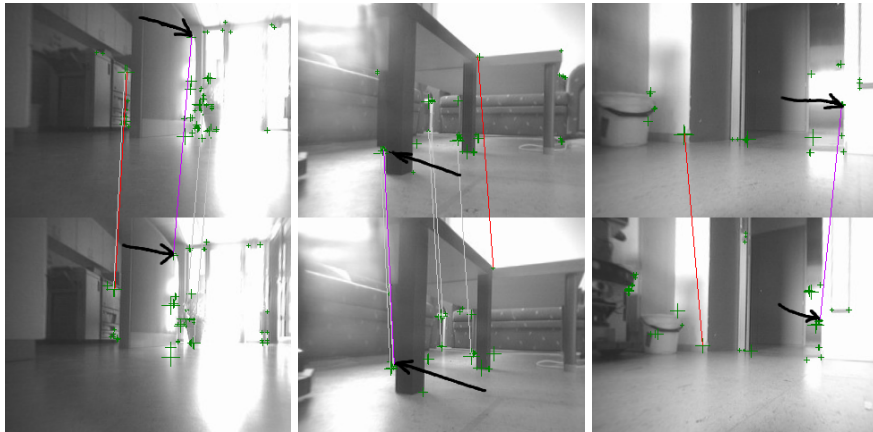


Fig. 3. Here the dark (red) and light (gray) lines show matches that were respectively fully initialized Cartesian points or used with the inverse depth parameterization. The purple lines (arrows) show features that have been identified as 'bad'.

add such constraint nodes for all the new correspondences that we found by SIFT matching and Lucas-Kanade tracking. Then we relax the pose states attached to these constraint nodes. The resulting energy of each constraint is then tested. If any are above some threshold we mark the image feature as bad and remove all its energy nodes. Image features marked as bad will be kept to be matched to again and any future image features that match to bad features will not be used. This will eliminate some of the occlusion type features, see figure 3.

Having tested the epipolar constraints we remove the temporary constraint nodes. We can then form two dimensional constraints between each pose and feature node. These are then stronger constraints than the epipolar one dimensional constraints between two pose nodes. They require forming 3D point feature state nodes. We add energy nodes to the 3D point feature nodes. These

energy nodes have an energy based on the innovation in the predicted pixel locations in the images. We do this without any initialization delay by using a special inverse depth representation as explained in the next section. We accumulate the information on each feature in this way. If the standard deviation in the estimated depth becomes less than 30% of the depth we initialize a full Cartesian represented point feature.

Figure 3 shows some examples of the constraints used to improve the odometry. In these example there were some initialized Cartesian points as well as a some inverse depth type constraints. There were also features labled as bad. These were image features that did not correspond to 3D points.

The benefit of features to the motion estimate decreases rapidly after about 3 or 4 matches are found between frames. On the other hand, matching to poor features or miss matching can cause the estimate to worsen. We therefore only keep the best matches (those that have the lowest constraint energy) plus any matches from any initialized Cartesian points. When features were scarce in the images we relaxed our threshold to allow some matches.

4 The Feature Representation

We did not begin by implementing the inverse depth parameterized features. Instead we began by implementing a conventional Cartesian coordinated representation for the features. This meant that we had to wait to initialize the point features until we had enough lateral motion to triangulate the feature location into good Gaussian ellipsoids in the Cartesian feature space. For features that did not resolve the depth sufficiently we simply left the epipolar constraint nodes on the graph. These pairwise constraints between poses are weaker than the constraints of forming a 3D feature but they do not require any depth. The results were an improvement over dead-reckoning alone but not nearly as good as that which we obtained by switching to the inverse depth parameterization.

The idea of the inverse depth parameterization is to use the pose the feature was first seen from as part of its parameterization. Then one takes the bearing angles or pixel locations in that frame as two other parameters. This eliminates the problem of the ever widening cone of possible feature location that a Cartesian xyz representation would have. By using the bearing angles as parameters the covariance is a well defined ellipse in the bearing space.

For the ninth and last parameter of the features one chooses the inverse of the depth of the feature. The reason for this is that the depth itself could initially be anything from a few cm to infinity. This range will be transformed to a finite range for the inverse depth. The feature can then be initialized at once without needing to accumulate information on it. The distribution in this parameter space is well approximated by a Gaussian from the start.

5 Experiments

The packbot robot can travel to inaccessible places, figure 1, up and down stairs, and on its back. Sensors are GPS, an IMU, and encoders on its two tracks. The GPS was not used. The encoders give reasonably good information on distance traveled while the robot moves straight on level ground with good traction. There is much error on turns. The IMU can reduce this but can not correct all of the errors. The robot is equipped with two cameras, one a wide angle camera and the other an IR camera. We used the IR camera.

We collected the images and telemetry data transmitted from the robot on a laptop computer connected via the wireless Ethernet. The images were 320X240 grey-scale jpg format. The first stage of the feature extraction was an undistortion of the images using the openCV library.

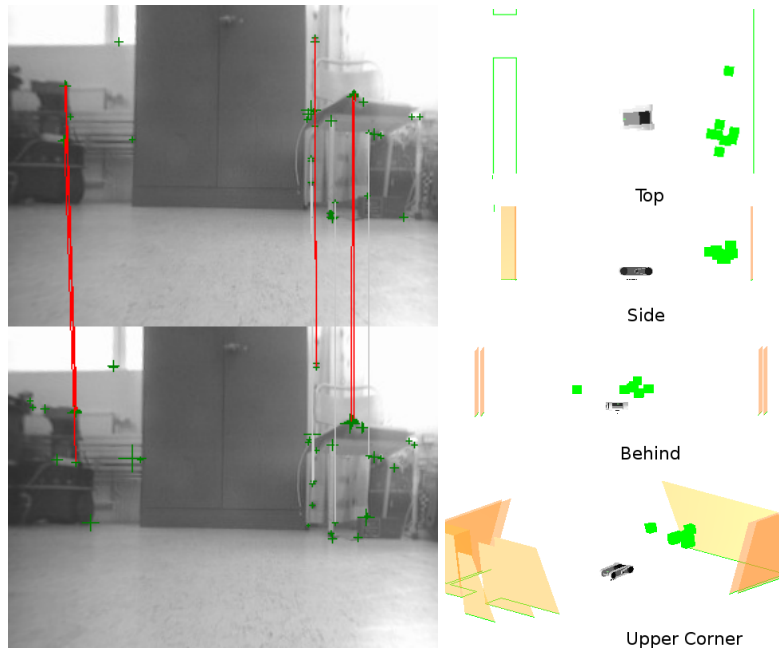


Fig. 4. Here we show how the matched image features get mapped in 3D. The early map is shown along with one of the early matches.

The first path was a simple circle around a table. By marking the floor we were able to return to exactly the same pose at the end of the test. This gave us a data set to work with and the ground truth of returning to the same pose. Note that global matching of the SIFT features was not done so the program would not recognize the features from the end as being the same as

those at the beginning. This global loop closing is future work. Our goal here is to try and minimize the drift of the estimate, a sort of visual odometry.

Figure 4 shows how the map looks in 3D just after the first few points have been added. The results are shown in figure 5. As you can see there were no major problems in this small test and the estimate is nearly perfect. In table 1 we summarize the error correction around the loop. The second test

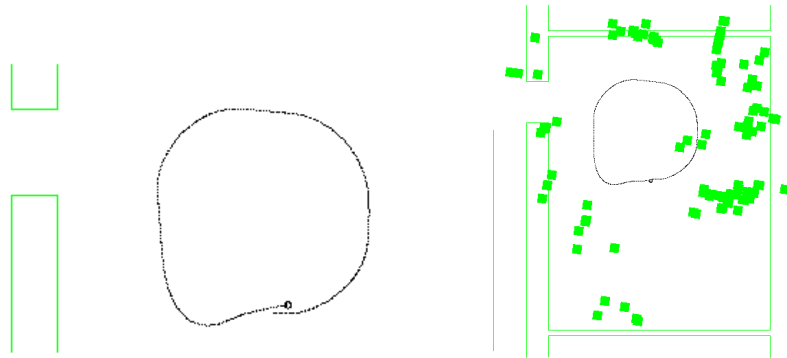


Fig. 5. Shown on the left is the dead-reckoning estimate of the robot path. The corrected path and the map of point features (shown as squares) are to the right.

Test (1230 frames)	x (m)	y (m)	z (m)	θ (rads)	ϕ (rads)	ψ (rads)
True $\pm < (1cm, .01rad)$	0.000	0.000	0.085	-1.4300	0.0000	0.0000
Dead-reckoning	0.104	-0.148	0.085	-1.3088	-0.0012	-0.0012
Inverse Depth Features	0.042	0.001	0.087	-1.4531	-0.0006	-0.0008

Table 1. This shows the estimated ending pose for the robot around the small loop. 117 entries were made in the global database, (initialized Cartesian points).

was a much longer path which lead out into the corridor and into two other rooms before returning to the exact starting location. The results are shown in figure 6. Here you see that the tuning of the filter had a significant impact on the results. A perfect loop closing without any global matching being done was not possible although we were able to get close on some occasions. The problem here as compared to the first case was not so much the length of the path but rather the variety of environments encountered. In some parts of the path there were simple no features found in the images. Not all occlusion type features could be filtered out of the data by the epipolar constraints.

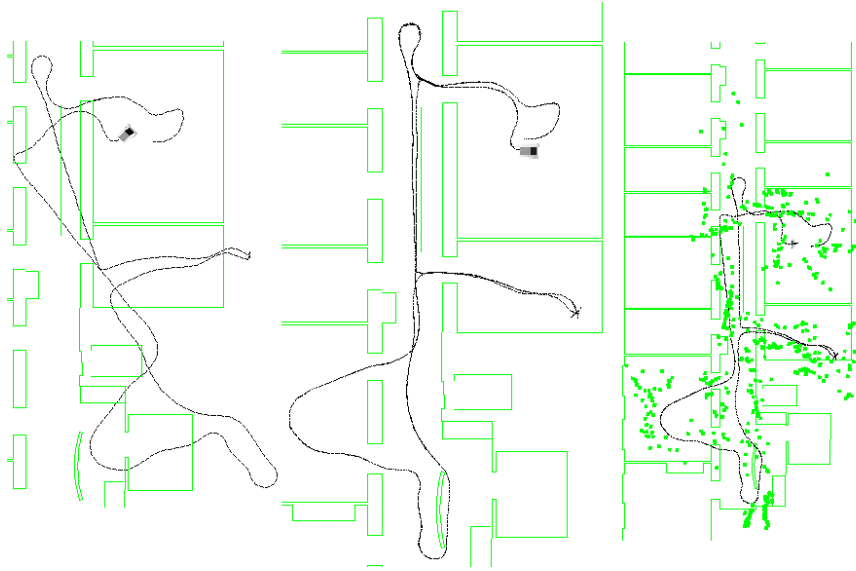


Fig. 6. Left we show the dead-reckoning. Center we show the best corrected path we were able to obtain using the image data. This map was the result of excessive parameter tuning but represents the ideal result we would like to get more robustly. On the bottom right we have a more typical good result of the corrected path that one obtains from using the image data, showing the features as well. Generally, one always gets some improvement regardless of the tuning but the amount varies from a little better to very good. The map on the right has been rotated 5 degrees about the start position in order to line the average path up with the building, better. These maps had about 650 fully initialized Cartesian features and 8,600 frames.

6 Conclusion

Despite some unavoidable short comings we saw a consistent improvement over dead-reckoning. With many good features the results were consistently good. This should give us a good base on which to build global matching and constraints on. The SIFT descriptors are after all designed specifically for global matching and should facilitate the next phase of this research.

The results indicated clearly that the inverse depth parameterization does solve the problem of initialization features seen in a single camera frame. This allowed us to replace one-dimensional epi-polar with two dimensional constraints. That in turn lead to better visual odometry estimates.

The use of Harris corners in combination with SIFT descriptors gave mixed results. We would like to have better criteria for selecting features in the images that would give features that were not only salient in the images but also gave bearings to true well defined 3D points. Our tracking method based on the constraint energy nodes eliminated most of the bad image features.

The fundamental problem we encountered was that of lack of any features at all in many of the images.

References

1. A. Davison, "Real-time simultaneous localization and mapping with a single camera," in *IEEE International Conference on Computer Vision*, 2003.
2. A. Davison, Y. G. Cid, and N. Kita, "Real-time 3D SLAM with a wide-angle vision," in *Intelligent Autonomous Vehicles (IAV-2004)*, M. I. Ribeiro and J. Santos-Victor, Eds., IFAC/EURON, IFAC/Elsevier, 2004.
3. D. Burschka and G. D. Hager, "V-GPS(SLAM): Vision-based inertial system for mobile robots," in *Proc. of the IEEE International Conference on Robotics and Automation (ICRA04)*, vol. 1, 2004, pp. 409–415.
4. P. Newman, D. Cole, and K. Ho, "Outdoor SLAM using visual appearance and laser ranging," in *Proc. of the IEEE International Conference on Robotics and Automation (ICRA06)*, vol. 1, 2006, pp. 1180–1188.
5. P. Jensfelt, D. Kragic, J. Folkesson, and M. Björkman, "A framework for vision based bearing only 3D SLAM," in *Proc. of the IEEE International Conference on Robotics and Automation (ICRA'06)*, Orlando, FL, May 2006.
6. J. Sola, A. Monin, M. Devy, and T. Lemaire, "Undelayed initialization in bearing only slam," in *Proc. of the IEEE/RJS International Conference on Intelligent Robots and Systems, (IROS 2005)*, vol. 1, 2005, pp. 2499–2504.
7. J. Montiel, J. Civera, and A. J. Davison, "Unified inverse depth parametrization for monocular SLAM," in *Robotics Science and Systems Conference*, 2006.
8. C. J. Harris and M. Stephens, "A combined corner and edge detector," in *Proc. 4th Alvey Vision Conference*, Manchester, 1988, pp. 147–151.
9. D. G. Lowe, "Object recognition from local scale-invariant features," in *Proc. of the Intl. Conf. on Computer Vision*, Corfu, Greece, 1999, pp. 1150–57.
10. S. Se, D. G. Lowe, and J. Little, "Mobile robot localization and mapping with uncertainty using scale-invariant visual landmarks," *International Journal of Robotics Research*, vol. 21, no. 8, pp. 735–58, 2002.
11. T. D. Barfoot, "Online visual motion estimation using fastslam with sift features," in *Proc. of the IEEE International Conference on Intelligent Robots and Systems (IROS2005)*, vol. 1. IEEE/RSJ, 2005, pp. 579–585.
12. R. Sim, R. Elinas, M. Griffin, and J. Little, "Vision-based slam using rao-blackwellised particle filter," in *Proc. of the IEEE International Conference on IJCAI Workshop Reasoning with Uncertainty*, vol. 1, 2005.
13. X. Wang and H. Zhang, "Good image features for bearing-only slam," in *Proc. of the IEEE International Conference on Intelligent Robots and Systems*, 2006.
14. J. Folkesson, P. Jensfelt, and H. I. Christensen, "Graphical SLAM using vision and the measurement subspace," in *Proc. of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS05)*, 2005.
15. J. Folkesson and H. I. Christensen, "Graphical SLAM for outdoor applications," *To Appear: Journal of Field Service Robotics*, 2006.
16. B. D. Lucas and T. Kanade, "An iterative image registration technique with and application to stereo vision," in *Proc. of the International Joint Conference on Artificial intelligence*, 1981.