

Self-Organizing Visual Maps

Robert Sim and Gregory Dudek

Centre for Intelligent Machines
McGill University
3480 University St, Montreal, Canada H3A 2A7

Abstract

This paper deals with automatically learning the spatial distribution of a set of measurements: images, in the examples presented here. The solution to this problem can be viewed as an instance of robot mapping although it can also be used in other contexts. We examine the problem of organizing an ensemble of images of an environment in terms of the positions from which the images were obtained, and where only limited prior odometric information is available. Our approach employs a feature-based method derived from a probabilistic robot localization framework. Initially, a set of visual landmarks are selected from the images and correspondences are found across the ensemble. The images are then localized by first assembling the small subset of images for which odometric confidence is high, and sequentially inserting the remaining images, localizing each against the previous estimates, and taking advantage of any priors that are available. We present experimental results validating the approach, and demonstrating metrically and topologically accurate results over two large image ensembles, even given only four initial ground truth poses. Finally, we discuss the results, their relationship to the autonomous exploration of an unknown environment, and their utility for robot localization and navigation.

1 Introduction

This paper addresses the problem of building a map of an unknown environment from an ensemble of observations and limited pose information. We examine the extent to which we can organize a set of measurements from an unknown environment to produce a visual map of that environment with little or no knowledge of where in the environment the measurements were obtained. In particular, we are interested in taking a set of snapshots of the environment using an uncalibrated monocular camera, and organizing them to quantitatively or qualitatively indicate where they were taken which, in turn, allows us to construct a visual map. We assume that, at most, we have limited prior trajectory information, so as to bootstrap the process—the source of this information might

be from the first few odometry readings along a trajectory, the general shape of the trajectory, information from an observer, or from a localization method that is expensive to operate, and hence is only applied at a small subset of the observation poses. While metric accuracy is of interest, our primary aim is to recover the topology of the ensemble. That is, metrically adjacent poses in the world are topologically adjacent in the resulting map.

The problem of automated robotic mapping is of substantial pragmatic interest for the development of mobile robot systems. The question of how we bootstrap a spatial representation, particularly a vision-based one, also appears to be relevant to other research areas such as computer vision and even ethology. Several authors have considered the use of self-organization in robot navigation [Takahashi *et al.*, 2001; Beni and Wang, 1991; Deneubourg *et al.*, 1989; Selfridge, 1962], often with impressive results. We believe this paper is among the first to demonstrate how to build a complete map of a real (non-simulated) unknown environment using monocular vision. We present quantitative data to substantiate this.

We approach the problem in the context of probabilistic robot localization using learned image-domain features (as opposed to features of the 3D environment) [Sim and Dudek, 2001]. To achieve this there are two steps involved: first, reliable features are selected and correspondences are found across the image ensemble. Subsequently, the quantitative behaviours of the features as functions of pose are exploited in order to compute a maximum-likelihood pose for each image in the ensemble. While other batch-oriented mapping approaches are iterative in nature [Thrun *et al.*, 1998; Kohonen, 1984], we demonstrate that if accurate pose information is provided for a small subset of images, the remaining images in the ensemble can be localized without the need for further iteration and, in some cases, without regard for the order in which the images are localized.

1.1 Outline

In the following section, we consider prior work related to our problem; in particular, approaches to self-organizing maps, and the simultaneous localization and mapping problem. We then proceed to present our approach, providing an overview of our landmark-based localization framework, followed by the details of how we apply the framework to organize the in-

put ensemble. Finally, we present experimental results on a variety of ensembles, demonstrating the accuracy and robustness of the approach.

2 Previous Work

The construction of self-organizing spatial maps (SOM's) has a substantial history in computer science. Kohonen developed a number of algorithms, for covering an input space [Kohonen, 1984; 1995]. While spatial coverage was used as a metaphor, the problem of representing a data space in terms of self-organizing features has numerous applications ranging from text searching to audition. The problem of spanning an input space with feature detectors or local basis functions has found wide application in machine learning, neural nets, and allied areas. In much of this algorithmic work, the key contributions have related to convergence and complexity issues.

The issue of automated mapping has also been addressed in the robotics community. One approach to fully automated robot mapping is to interleave the map synthesis and position estimation phases of robot navigation (sometimes known as SLAM: simultaneous localization and mapping). As it is generally applied, this entails incrementally building a map based on geometric measurements (e.g. from a laser rangefinder, sonar or stereo camera) and intermittently using the map to correct the robot's position as it moves [Leonard and Durrant-Whyte, 1991b; Yamauchi *et al.*, 1998; Davison and Kita, 2001]. When the motion of a robot can only be roughly estimated, a topological representation becomes very attractive. Early work by Kuipers and Byun used repeated observation of a previously observed landmark to instantiate cycles in a topological map of an environment during the mapping process [Kuipers and Byun, 1987; 1991]. The idea of performing SLAM in a topological context was also been examined theoretically [Deng and Mirzaian, 1996; Dudek and Tsotsos, 1991]. The probabilistic fusion of uncertain motion estimates has been examined by several authors (cf. [Smith and Cheeseman, 1986]) and the use of Expectation Maximization has recently proven quite successful although it still depends on estimates of successive robot motions [Shatkay and Kaelbling, 1997; Thrun, 1998; Choset and Nagatani, 2001].

3 Landmark Framework

Our approach employs an adaptation of the landmark learning framework described in [Sim and Dudek, 2001]. We review it here in brief and refer the reader to the cited work for further details.

The key idea is to learn visual features, parametrically describe them so that they can be used to estimate one's position (that is, they can be used for localization). The features are pre-screened using an attention operator that efficiently detects statistically anomalous parts of an image and robust, useful features are recorded along with an estimate of their individual utility.

In the localization context, assume for the moment that we have collected an ensemble of training images with ground-truth position information associated with each image. In

other words, we have the full solution to the mapping problem in hand (we will relax this assumption in the following sections). The landmark learning framework operates by first selecting a set of local features from the images using a measure of visual attention, tracking those features across the ensemble of images by maximizing the correlation of the local image intensity of the feature, and subsequently parameterizing the set of observed features in terms of their behaviour as a function of the known positions of the robot (Figure 1).

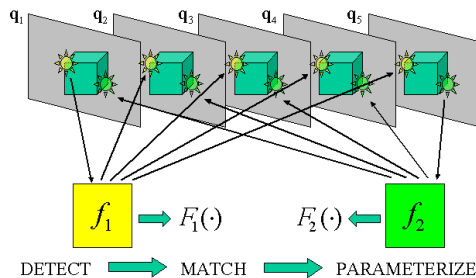


Figure 1: Landmark Learning Framework.

The resulting *tracked landmarks* can be applied in a Bayesian framework to solve the localization problem. Specifically, given an observation image \mathbf{z} , the probability that the robot is at pose \mathbf{q} is proportional to the probability of the observation conditioned on the pose:

$$p(\mathbf{q}|\mathbf{z}) = \frac{p(\mathbf{z}|\mathbf{q})p(\mathbf{q})}{p(\mathbf{z})} \quad (1)$$

where $p(\mathbf{q})$ is the prior on \mathbf{q} and $p(\mathbf{z})$ is a normalization constant. For a landmark-based approach, we express the probability of the observation conditioned on the pose as a mixture model of probability distributions derived from the individual landmarks:

$$p(\mathbf{z}|\mathbf{q}) = k \sum_{\mathbf{l}_i \in \mathbf{z}} p(\mathbf{l}_i|\mathbf{q}) \quad (2)$$

where \mathbf{l}_i is a detected observation of landmark i in the image and k is a normalizing constant.

The individual landmark models are generative in nature. That is, given the proposed pose \mathbf{q} , an expected observation \mathbf{l}_i^* is generated by learning a parameterization $\mathbf{l}_i^* = F_i(\mathbf{q})$ of the landmark, and the observation probability is determined by a Gaussian distribution centered at the expected observation and with covariance determined by cross-validation over the training observations. Whereas in prior work the parameterization was computed using radial basis function networks, in this work we construct the interpolants using a Delaunay triangulation of the observation poses. In general, the contents of the feature vector \mathbf{l}_i itself can be any quantitative measure related to the local image neighbourhood of the observation. For the present work we define the feature vector as the position of the feature in the image:

$$\mathbf{l}_i = [x_i \ y_i] \quad (3)$$

A pose estimate is obtained by finding the pose \mathbf{q}^* that maximizes Equation 1. It should be noted that the framework

requires no commitment as to how uncertainty is represented or the optimization is performed. It should be noted, however, that the probability density for \mathbf{q} might be multi-modal, and, as is the case for the problem at hand, weak priors on \mathbf{q} might require a global search for the correct pose. For this work, we employ a multi-resolution grid decomposition of the environment, first approximating $p(\mathbf{q}|\mathbf{z})$ at a coarse scale and computing increasingly higher resolution grids in the neighbourhood of \mathbf{q}^* as it is determined at each resolution. Figure 2 illustrates a typical probability distribution over the pose space, given an input image.

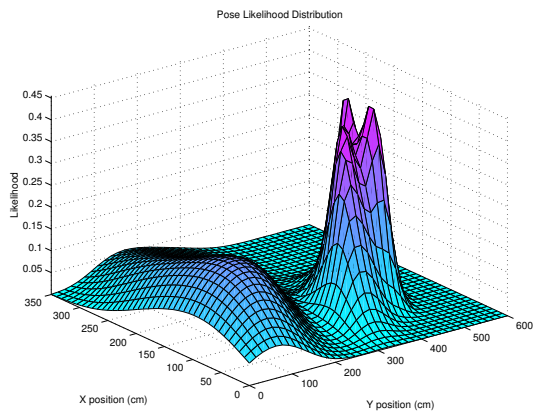


Figure 2: Likelihood function of pose of robot over 3.0m by 6.0m pose space. Note that the distribution is not unimodal.

4 Self Organization

We now turn to the problem of inferring the poses of the training images when ground truth is unavailable, or only partially available. The self-organization process involves two steps. In the first step, image features are selected and tracked, and in the second step the set of images are localized.

4.1 Tracking

Tracking proceeds by considering the images in an arbitrary order (possibly, but not necessarily, according to distance along the robot’s trajectory). An attention operator is applied to the first image \mathbf{z} in the set¹, and each detected feature initializes a *tracking set* $T_i \in T$. The image itself is added to the ensemble set E . For each subsequent image \mathbf{z} , the following algorithm is performed:

1. A search is conducted over the image for matches to each tracking set in T , and successful matches are added to their respective tracking sets T_i . Call the set of successful matches M .
2. The attention operator is then applied to the image and the set of detected features S is determined.
3. If the cardinality of M is less than the cardinality of S , new tracked sets T_i are initialized by elements selected from S . The elements are selected first on the basis of

¹We select local-maxima of edge-density.

their response to the attention operator, and second on the basis of their distance from the nearest image position in M . In this way, features in S which are close to prior matches are omitted, and regions of the image where features exist but matching failed receive continued attention. Call this new set of tracking sets T_S .

4. A search for matches to the new tracking sets in T_S is conducted over each image in E (that is, the previously examined images), and the successful matches are added to their respective tracking set.
5. $T = T \cup T_S$
6. $E = E \cup \mathbf{z}$

The template used for by any particular tracking set is defined as the local appearance image of the initial feature in the set. We use local windows of 33 pixels in width and height. Matching is considered successful when the normalized correlation of the template with the local image under consideration exceeds a user-defined threshold.

When tracking is completed, we have a set of feature correspondences across the ensemble of images. The process is $O(kn)$ where k is the final number of tracked sets, and n is the number of images.

4.2 Localization

Once tracking is complete, the next step is to determine the position of each image in the ensemble. For the moment, consider the problem when there is a single feature that was tracked reliably across all of the images. If we assume that the image feature is derived from a fixed 3D point in space, the motion of the feature through the image will be according to a monotonic mapping as a function of camera pose and the camera’s intrinsic parameters. As such, the topology of a set of observation poses is preserved in the mapping from pose-space to image-space. While the mapping itself is nonlinear (due to perspective projection), it can be approximated by associating actual poses with a small set of the observations and determining the local mappings of the remaining unknown poses by constructing an interpolant over the known poses. Such an algorithm would proceed as follows:

1. Initialize $S = \{(\mathbf{q}, \mathbf{z})\}$, the set of (pose, observation) pairs for which the pose is known. Compute D , the parameterization of S as defined by the landmark learning framework.
2. For each observation \mathbf{z} with unknown pose,
 - (a) Use D as an interpolant to find the pose \mathbf{q}^* that maximizes the probability that \mathbf{q}^* produces observation \mathbf{z} .
 - (b) Add $(\mathbf{q}^*, \mathbf{z})$ to S and update D accordingly.

For a parameterization model based on a Delaunay Triangulation interpolant, updating the D takes $O(\log n)$ amortized time, where n is the number of observations in the model. The cost of updating the covariance associated with each model is $O(k \log n)$, where k is the number of samples omitted during cross-validation.



Figure 3: The first scene, a 2.0m by 2.0m pose space.

In addition, the cost of finding the maximally likely pose with D is $O(m \log n)$, where m corresponds to the number of poses that are evaluated (finding a face in the triangulation that contains a point \mathbf{q} can be performed in $\log n$ time.). Given n total observations, the entire algorithm takes $O(n(m + k + 1) \log n)$ time. Both m and k can be bounded by constants, although in practice we typically bound k by n .

In practice, of course, there is more than one feature detected in the image ensemble. Furthermore, in a suitably small environment, some might span the whole set of images, but in most environments, most are only visible in a subset of images. Finally, matching failures might introduce a significant number of outliers to individual tracking sets. Multiple landmarks, and the presence of outlier observations are addressed by the localization framework we have presented; the maximum likelihood pose is computed by maximizing Equation 2, and the effects of outliers in a tracked set are reduced by their contribution to the covariance associated with that set.

When it cannot be assumed that the environment is small enough such that one or more landmark spans it, we must rely on stronger priors to bootstrap the process. For example, we might require the initial known poses to be close together, ensuring that they share common landmarks for parameterization. In addition, we might take advantage of knowledge of the order in which images were acquired along a trajectory, ensuring that as one landmark goes out of view, new ones are present against which to localize.

In the following section we present experimental results on two image ensembles.

5 Experimental Results

5.1 A Small Scene

For our first experiment, we demonstrate the procedure on a relatively compact scene. An ensemble of 121 images of the scene depicted in Figure 3 was collected over a 2m by 2m environment, at 20cm intervals. Ground truth was measured by hand, accurate to 0.5cm.

Given the ensemble, the images were sorted at random and tracking was performed as described in Section 4.1, resulting in 91 useful tracked landmarks. (A tracked landmark

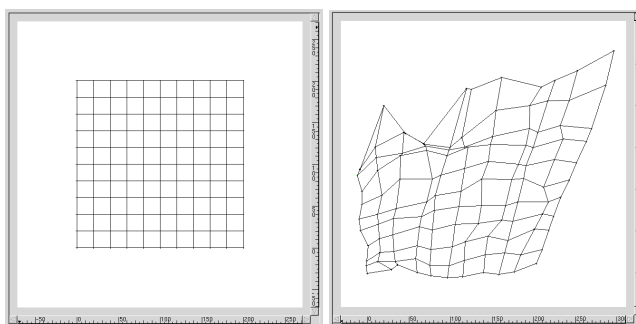


Figure 4: Ground truth, and the map resulting from the self-organizing process for the environment depicted in Figure 3.



Figure 5: The second scene, a 3.0m by 5.5m pose space.

was considered useful if it contained at least 4 observations). The localization stage proceeded by first providing the ground truth information to four images selected at random. The remaining images were again sorted at random and added, without any prior information about their pose, according to the methodology described in the previous section. Figure 4 plots the original grid of poses, and beside it the same grid imposed upon the set of pose estimates computed for the ensemble. While there is some warping in the mesh, for the most part the topology of the poses is clearly preserved. It is interesting to note that the mesh is distorted most as the y-axis increases, corresponding to looming forward with the camera and, as such, where the nonlinearity of the perspective projection is most pronounced.

5.2 A Larger Scene

For our second experiment, we examine a larger pose space, 3.0m in width and 5.5m in depth, depicted in Figure 5. For this experiment, 252 images were collected at 25cm intervals using a pair of robots, one of which used a laser range-finder to measure the 'ground-truth' pose of the moving robot [Rekleitis *et al.*, 2001].

As in the previous experiment, tracking was performed over the image ensemble and a set of 49 useful tracked landmarks were extracted. In this instance, the larger interval between images, some illumination variation in the scene and the larger number of input images presented significant

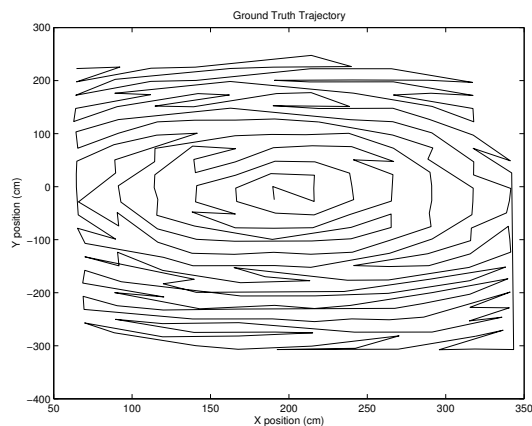


Figure 6: Ground truth simulated trajectory.



Figure 7: Evolution of the annular prior over the first few input images.

challenges for the tracker, resulting in the smaller number of tracked landmarks.

Given the size of the environment, no one landmark spanned the entire pose space. As a result, it was necessary to impose constraints on the input ground-truth priors, and the order in which the input images were considered for localization. In addition, a weak prior $p(\mathbf{q})$ was applied as each image was added in order to control the distortion in the mesh.

Rather than select the initial ground-truth images at random, ground truth was supplied to the four images closest to the centre of the environment. The remainder of the images were sorted by simulating a spiral trajectory of the robot through the environment, intersecting each image pose, and adding the images as they were encountered along the trajectory. Figure 6 illustrates the simulated ground-truth trajectory through the ensemble. Finally, given the sort order, as images were added it was assumed that their pose fell on an annular ring surrounding the previously estimated poses. The radius and width of the ring was defined in terms of the interval used to collect the images. The computed priors over the first few images input into the map are depicted in Figure 7. The intent of using these priors was to simulate a robot exploring the environment along trajectories of increasing radius from a home position.

As in the previous section, Figure 8 plots the original grid of poses, and beside it the same grid imposed upon the set of pose estimates computed for the ensemble. Again, the positive y -axis corresponds to looming forward in the image, and as such the mesh distorts as landmarks accelerate in image space as the camera approaches them. Note however, that as in the first experiment, the topology of the poses preserved for most of the grid.

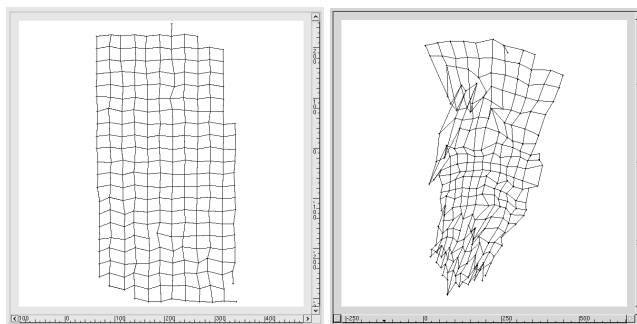


Figure 8: Ground truth, and the map resulting from the self-organizing process for the environment depicted in Figure 3.

5.3 Discussion

We have demonstrated an approach to spatially organizing images of an unknown environment using little or no positional prior knowledge. The repeated occurrences of learned visual features in the images allows us to accomplish this. The visual map of the environment that is produced appears to be topologically correct and also demonstrates a substantial degree of metric accuracy and can be described as a locally conformal mapping of the environment. This representation can then be readily used for path execution, trajectory planning and other spatial tasks.

While several authors have considered systems that interleave mapping and position estimation, we believe ours is among the first to do this based on monocular image data. In addition, unlike prior work which typically uses odometry to constrain the localization process, we can accomplish this with essentially no prior estimate of the position the measurements are collected from. On the other hand, if some positional prior is available we can readily exploit it. In the second example shown in this paper we exploited such a prior. Even in this example, it should be noted that the data acquisition trajectory was one that did not include cycles. In general, cyclic trajectories (ones that re-visit previously seen locations via another route) will greatly improve the quality of the results; in fact they are prerequisite for many existing mapping and localization techniques, both topological and metric ones.

We believe that absence of a requirement of a position prior (i.e. odometry) makes this approach suitable for unconventional mapping applications, such as the integration of data from walking robots or from manually collected video sequences. Our ability to do this depends on the repeated occurrence of visual features in images from adjacent positions. This implies that successful mapping depends on images being taken at sufficiently small intervals to assure common elements between successive measurements.

References

- [Beni and Wang, 1991] G. Beni and J. Wang. Theoretical problems for the realization of distributed robotic system. In *Proc. IEEE Int'l Conf on Robotics and Automation*, pages 1914–1920, Sacramento, CA, April 1991. IEEE Press.

- [Choset and Nagatani, 2001] H. Choset and K. Nagatani. Topological simultaneous localization and mapping (slam): toward exact localization without explicit localization. *IEEE Tr. on Robotics and Automation*, 17(2), April 2001.
- [Davison and Kita, 2001] A. J. Davison and N. Kita. 3D simultaneous localisation and map-building using active vision for a robot moving on undulating terrain. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2001.
- [Deneubourg *et al.*, 1989] J.L. Deneubourg, S. Goss, J. Pasteels, D. Fresneau, and J.P. Lachaud. Self-organization mechanisms in ant societies (ii): Learning in foraging and division of labor. In J.M. Pasteels and eds. J.L. Deneubourg, editors, *From Individual to Collective Behavior in Social Insects, Experientia Supplementum*, volume 54, pages 177–196. Birkhauser Verlag, 1989.
- [Deng and Mirzaian, 1996] X. Deng and A. Mirzaian. Competitive robot mapping with homogeneous markers. *IEEE Trans. on Robotics and Automation*, 12(4):532–542, August 1996.
- [Dudek and Tsotsos, 1991] Gregory Dudek and John K. Tsotsos. Shape representation and recognition from curvature. In *Proceedings of the 1991 Conference on Computer Vision and Pattern Recognition*, pages 35–41, Maui, Hawaii, June 1991. IEEE.
- [Gross *et al.*, 1996] H. Gross, V. Stephan, and H. Bohme. Sensory-based robot navigation using selforganizing networks and q-learning. 1996.
- [Kohonen, 1984] T. Kohonen. *Self-Organization and Associative Memory*. Springer-Verlag, New York, 1984.
- [Kohonen, 1995] Teuvo Kohonen. *Self-organizing maps*. Springer, Berlin; Heidelberg; New-York, 1995.
- [Kuipers and Beeson, 2002] Benjamin Kuipers and Patrick Beeson. Bootstrap learning for place recognition. In *Proc. National Conference on Artificial Intelligence (AAAI)*, 2002.
- [Kuipers and Byun, 1987] Benjamin J. Kuipers and Y. T. Byun. A qualitative approach to robot exploration and map-learning. In *Proceedings of the IEEE workshop on spatial reasoning and multi-sensor fusion*, pages 390–404, Los Altos, CA, 1987. IEEE.
- [Kuipers and Byun, 1991] B. Kuipers and Y-T. Byun. A robot exploration and mapping strategy based on a semantic hierarchy of spatial representations. *Robotics and Autonomous Systems*, 8:47–63, 1991.
- [Leonard and Durrant-Whyte, 1991a] J. Leonard and H. F. Durrant-Whyte. Simultaneous map building and localization for an autonomous mobile robot. In *Proc. IEEE Int. Workshop on Intelligent Robots and Systems*, pages 1442–1447, Osaka, Japan, November 1991.
- [Leonard and Durrant-Whyte, 1991b] John J. Leonard and Hugh F. Durrant-Whyte. Mobile robot localization by tracking geometric beacons. *IEEE Transactions on Robotics and Automation*, 7(3):376–382, June 1991.
- [Leonard and Feder, 2000] John J. Leonard and Hans Jacob S. Feder. A computationally efficient method for large-scale concurrent mapping and localization. In John Hollerbach and Dan Koditschek, editors, *Robotics Research: The Ninth International Symposium*, London, 2000. Springer-Verlag.
- [Rekleitis *et al.*, 2001] Ioannis Rekleitis, Robert Sim, Gregory Dudek, and Evangelos Milios. Collaborative exploration for the construction of visual maps. In *2001 IEEE/RSJ Conference on Intelligent Robots and Systems (IROS)*, Hawaii, October 2001.
- [Selfridge, 1962] O.G. Selfridge. The organization of organization. In M.C. Yovits, G.T. Jacobi, and G.D. Goldstein, editors, *Self-Organizing Systems*, pages 1–7. McGregor & Werner, Washington D.C., 1962.
- [Shatkay and Kaelbling, 1997] Hagit Shatkay and Leslie Pack Kaelbling. Learning topological maps with weak local odometric information. In *Proc Int. Joint Conf. Artificial Intelligence (IJCAI)*, pages 920–929, Nagoya, Japan, 1997. Morgan Kaufmann.
- [Sim and Dudek, 2001] R. Sim and G. Dudek. Learning generative models of scene features. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, Lihue, HI, December 2001. IEEE Press.
- [Smith and Cheeseman, 1986] Randall C. Smith and Peter Cheeseman. On the representation and estimation of spatial uncertainty. *International Journal of Robotics Research*, 5(4):56–68, 1986.
- [Takahashi *et al.*, 2001] Takashi Takahashi, Toshio Tanaka, Kenji Nishida, and Takio Kurita. Self-organization of place cells and reward-based navigation for a mobile robot. *Proceedings of the 8th International Conference on Neural Information Processing, Shanghai (China)*, November 14–18 2001.
- [Thrun *et al.*, 1998] S. Thrun, D. Fox, and W. Burghard. A probabilistic approach to concurrent mapping and localization for mobile robots. *Autonomous Robots*, 5:253–271, 1998.
- [Thrun, 1998] S. Thrun. Learning metric-topological maps for indoor mobile robot navigation. *AI Journal*, 99(1):21–71, 1998.
- [Yamauchi *et al.*, 1998] B. Yamauchi, A. Schultz, and W. Adams. Mobile robot exploration and map building with continuous localization. In *Proc. IEEE Int. Conf. on Robotics and Automation*, pages 3715–2720, Leuven, Belgium, May 1998. IEEE Press.