

Low-Level Visual Homing

Andrew Vardy and Franz Oppacher

School of Computer Science
Carleton University
Ottawa, Canada

avardy@scs.carleton.ca
<http://www.scs.carleton.ca/~avardy>

Abstract. We present a variant of the *snapshot model* [1] for insect visual homing. In this model a snapshot image is taken by an agent at the goal position. The disparity between current and snapshot images is subsequently used to guide the agent's return. A matrix of local low-level processing elements is applied here to compute this disparity and transform it into a motion vector. This scheme contrasts with other variants of the snapshot model which operate on one-dimensional images, generally taken as views from a synthetic or simplified real world setting. Our approach operates directly on two-dimensional images of the real world. Although this system is not a model of any known neural structure, it hopes to offer more biological plausibility than competing techniques because the processing applied is low-level, and because the information processed appears to be of the same sort of information that is processed by insects. We present a comparison of results obtained on a set of real-world images.

1 Introduction

In [1] the *snapshot model* was proposed to explain the remarkable ability of honeybees to return to a place of interest such as a nest or food source after being displaced from that place. Variants of this model have been implemented both in simulation and in physical robots [4,11,3,7,14,8,9,10]. All of these computational models operate on one-dimensional images. Even for the case of physical robots with cameras that generate two-dimensional images, the height dimension is almost immediately collapsed such that the image used for homing is effectively one-dimensional. Our approach differs from these in that we apply our processing directly on two-dimensional images. We find that the use of two-dimensional image features improves results over a similar model that operates only on one-dimensional features. Also, a number of these other models are tested either in simulated worlds or in simplified real world scenarios. We test our approach using the real world image album of [10] and compare our results with theirs.

The snapshot model was developed to match data of honeybee search patterns. A model agent is placed at the goal and allowed to capture a snapshot image. It is then displaced and allowed a return attempt. The model operates on one-dimensional panoramic binary images. Features (region centers) in the

current image are paired with their closest matching counterparts in the snapshot image. Each pairing generates two vectors: one correcting for bearing and one correcting for differences in apparent size. The sum of all vectors yields the motion vector. One key requirement of the snapshot model is that the agent maintains a consistent orientation. There is evidence to suggest that bees take on the same orientation when returning to a remembered place as they took when originally learning the layout of that place [2,15,6]. A robot homing via the snapshot model must employ some sort of compass system to maintain, or compensate for changes, in orientation.

In [8] a neural implementation of the snapshot model was presented. It showed that the processing necessary to implement the snapshot model could be achieved using simple neuron-like elements arranged in layers of interconnected rings. The outermost ring was exposed to a one-dimensional panoramic image of the homing agent's environment. Subsequent layers processed this image and extracted the necessary vectors to achieve homing. Our approach has been directly inspired by this work. It differs, however, in its use of two-dimensional images, and hence, two-dimensional processing layers. Also, the use of real-world test images prompted a number of modifications.

We can not state that our method is a model of any known neural circuit in an insect's brain. According to [9], "...*nothing is known so far about the neural circuits that realize visual homing in insect brains.*" Still, there are three ways in which we claim that our model has some measure of biological plausibility. First, all computations are made via a matrix of locally-connected neuron-like elements. Second, this matrix has a layered structure which preserves spatial relationships from layer to layer. This is inspired by the retinotopic structure of the mammalian visual cortex where neurons on one layer seem to share receptive fields with the layers below them [5]. Our final claim to biological plausibility is in regard to the sort of information that our method employs. It operates on snapshots of a scene, stored in retinal coordinates. This is the primary contention of the snapshot model [1]. More recent work on wood ants [6] suggests the storage and use of two-dimensional templates for ant homing. We hope that our model can help to close the explanatory gap between a complete neural map of the hardware that implements visual homing in insects and a purely computational model that cares nothing for the details of implementation.

Our method operates by approximating, for each feature in the snapshot image, the direction in which that feature has moved in the current image. The features we use are image corners. The method used to determine the direction of feature movement is to examine the local gradient, created from features in the current image, at the position of the snapshot feature. Thus, we refer to our method here as the *Corner Gradient Snapshot Model*, or CGSM.

2 Processing Matrix

We will now describe the operation of the processing matrix for CGSM. The processing that is applied is depicted in figure 1. A summary of this processing

follows: An input image is fed into the processing matrix. It is first smoothed and then corners are extracted as distinct features. If the agent is at the goal then the image of features is stored as the snapshot image. Gradients are formed around each feature and these gradients are locally compared with features in the snapshot image to generate vectors which indicate the direction that these features have moved in. These vectors specifying motion in the image are mapped onto vectors specifying the corresponding motion of the agent. Finally, this last set of vectors is summed to create the agent's overall motion vector. For the following layer descriptions, the image that is fed in from the previous layer is referred to as the *input image*. There is insufficient space to describe exactly how each layer's functionality can be distributed across a grid of low-level processing elements. We hope that the simplicity of each layer will provide the reader with assurance that this distribution is possible.

Gaussian Filter. Convolves the input image with a 5x5 Gaussian mask.

Corner Extraction. Extracts corners from the input image. The corner detection scheme we use is known as the Harris detector and is essentially the same as that presented in [13] (pages 82-85). The required computation is purely local in scope and requires only the application of some simple hard-coded algebra to the results of local convolutions and sums.

Local Maxima Extraction. Applies thresholding and non-maxima suppression to find peaks in the input image [13].

Gradient Formation. Forms decaying gradients around points (maxima) in the input image. These gradients allow the next layer to detect the direction in which a corner has moved from its idealized position at a stored snapshot point. Repeated gaussian filtering and diffusion are two possible ways of implementing this layer. We opt for a different approach which is quicker to compute with a serial computer. A single gradient mask is constructed with a center value of unity that decays out radially according to $1/d^\zeta$, where d is the distance from center and ζ is an arbitrary constant ($0 > \zeta > 1$). Beginning with a zero image, copies of the mask are centered over each feature. The value of each pixel in the output image is taken as the maximum value from all mask copies aligned directly above it. Note that we assume that the agent moves only horizontally in the plane. Thus, corners in the top half of the image should never be paired with those in the bottom half. We implicitly prevent such pairings by forming gradients in the top and bottom halves of the image separately.

Ring Operator. Applied at every non-zero point in the snapshot image to generate a vector pointing in the direction that the feature has moved. This layer actually has two input images: the snapshot image, and the gradient image. At each non-zero point in the snapshot image, S , a ring of detector cells surrounding S is employed to detect the approximate direction of the gradient at that point. Ideally, this direction will point to the matching feature F . The vector from S to the detector cell with the highest response yields the uphill direction of the gradient. The vectors generated by this layer are *image vectors* in that they describe the motion of features within

the plane of the image. Image vectors will be indicated below in lower-case (\vec{u} , \vec{v} , \vec{i} , and \vec{e}).

Vector Mapping. Maps image vectors into *agent motion vectors*. An agent motion vector describes the motion of the agent within its plane of travel that would correct (reduce to zero) the corresponding image vector. The next section provides more detail on this layer. Agent motion vectors will be indicated in upper-case (\vec{V} , \vec{T} , and \vec{E}).

Vector Sum. Sums the agent motion vectors in the input image to generate a single vector that will move the agent toward home.

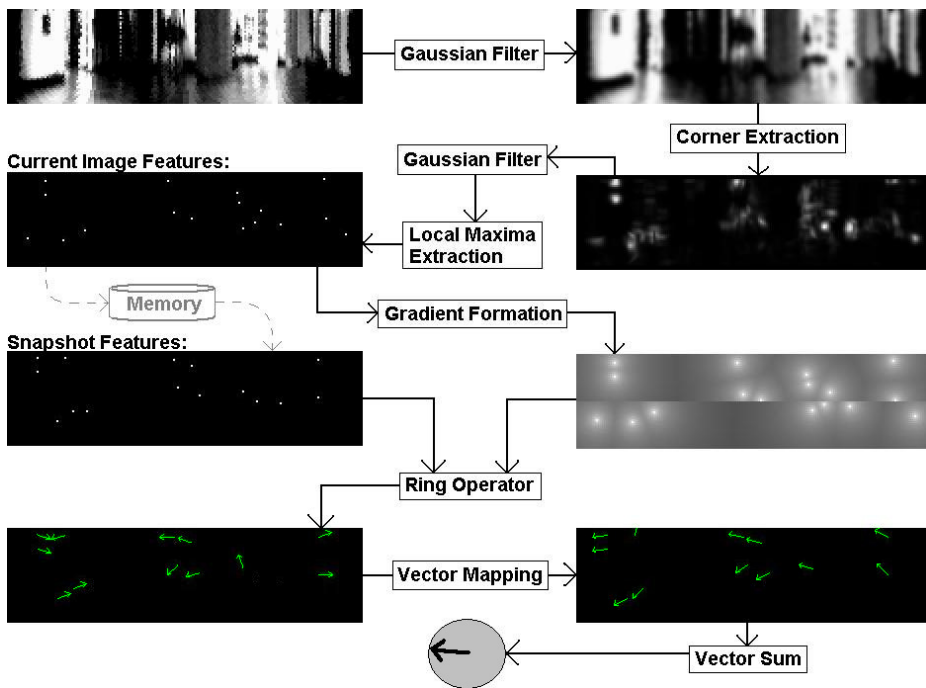


Fig. 1. Processing applied by the CGSM matrix. The input is the image in the upper-left corner. The output is the movement vector in the shaded circle at the bottom. Processing layers are shown as boxed text. Dashed arrows relate to the storage and recall of the snapshot image. The input image was taken from position (270,180) in the image album discussed in the *Results* section. This image was taken from a point directly to the right of the snapshot position (210, 180). Hence, the direction of the output movement vector is approximately correct.

2.1 Vector Mapping

Each image vector \vec{u} is attached or associated with the position of a feature S in the snapshot image. $\vec{v} = -\vec{u}$ is a unit vector which points from C to S , where C is the position of a feature in the current image. If features are paired correctly then S and C will both correspond to a single *environmental feature* F . In this section we will assume that this is true. The vector mapping layer computes an agent motion vector \vec{V} , using \vec{v} and S_x (horizontal position of snapshot feature within the image). Ideally each \vec{V} would point directly home. However, there is insufficient information (e.g. distance to F) to compute \vec{V} exactly, therefore we develop an approximation. First we decompose \vec{v} and \vec{V} .

$$\vec{v} = a\vec{i} + b\vec{e} \tag{1}$$

$$\vec{V} = \alpha\vec{I} + \beta\vec{E} \tag{2}$$

\vec{I} is a movement of the agent towards a feature X (the relationship of X to S and C will be made clear below). With this movement, X will move upwards in the image according to \vec{i} . One can imagine walking towards a light and having the image of that light travel straight upwards in one’s field of vision. The walking movement is \vec{I} and the corresponding movement of the light is \vec{i} . \vec{E} is a movement of the agent 90° to the right of X . With this movement, X will move to the left according to \vec{e} . We can similarly imagine walking 90° to the right of a light and seeing the image of that light travel to the left in our field of vision. In this case, the walking movement is \vec{E} and the movement of the light is \vec{e} . Figure 2 illustrates these relationships. The quantity θ_X is the horizontal angular position of X in the image, $\theta_X = \frac{2\pi X_x}{w}$, where X_x is the x -coordinate of X and w is the image width. We can now give the form of \vec{i} , \vec{e} , \vec{I} , and \vec{E} .

$$\begin{aligned} \vec{I}(\theta_X) &= \begin{bmatrix} \sin \theta_X \\ \cos \theta_X \end{bmatrix} & \vec{i} &= \begin{bmatrix} 0 \\ 1 \end{bmatrix} \\ \vec{E}(\theta_X) &= \begin{bmatrix} \sin (\theta_X + 90^\circ) \\ \cos (\theta_X + 90^\circ) \end{bmatrix} & \vec{e} &= \begin{bmatrix} -1 \\ 0 \end{bmatrix} \end{aligned}$$

We wish to have the feature C move in the image toward S . However, the exact position of C is unknown. We make the assumption that C is not too distant from S . That is, $S \approx C$. Thus, any movement of the agent \vec{V} that would make S move along \vec{v} would also make C move along \vec{v} . Of course, S cannot move at all because it is a feature recalled from memory and is fixed in the image. The feature X is a stand-in for S intended to make it less awkward to speak about the movement of S . Summarizing: $X = S$ and $X \approx C$.

We now return to find the constants a , b , α , and β in equations 1 and 2. Given \vec{i} and \vec{e} we can solve equation 1 for a and b . The solution: $a = v_y$ and $b = -v_x$. We make a further approximation to find α and β . If a is much larger than b than the difference in vertical image position between S and C is greater than the difference in horizontal position. In this case, the agent’s distance to the environmental feature F along \vec{I} will be greater than the distance along

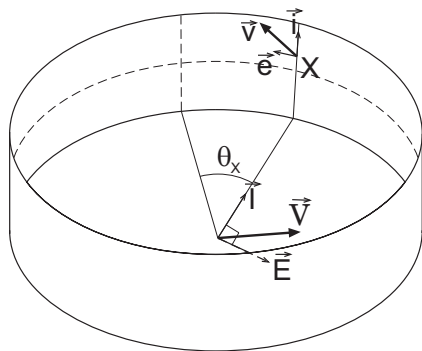


Fig. 2. Motion of the image feature X in direction \vec{v} is generated by motion of the agent in direction \vec{V} .

\vec{E} . The opposite will be true if b is much larger than a . Thus, we make the approximation: $\alpha = a$ and $\beta = b$. The final mapping from \vec{v} to \vec{V} is as follows.

$$\vec{V} = \begin{bmatrix} v_y \sin \theta_S - v_x \sin (\theta_S + 90^\circ) \\ v_y \cos \theta_S - v_x \cos (\theta_S + 90^\circ) \end{bmatrix} \quad (3)$$

One last caveat remains. The vector \vec{i} points upwards in figure 2 because movement in direction \vec{V} would cause X to sweep upwards. However, if X had been a feature in the bottom half of the image then X would have swept downwards. Thus, for features in the bottom half of the image \vec{i} will be $(0, -1)$.

3 Results

We compare the homing performance of CGSM against an implementation of the average landmark vector model which was designed to operate on real-world images [10,12]. The average landmark vector model can be considered a variant of the snapshot model [9]. This approach extracts features from a 1-D image that is itself extracted from a 2-D panoramic image. Once the features have been extracted the average landmark vector model is applied in its original form [7]. We shall refer to this implementation of the average landmark vector model as *XALV* (eXtracted Average Landmark Vector).

The authors of [10] have very kindly made their images available to us and we use them here to compare our results with theirs. These images were taken from a grid of positions spaced 30 cm apart within a $9m \times 3m$ area. The environment was an unmodified entrance hall of a university building. These pictures were taken by a robot-mounted camera pointed upwards at a conical mirror. The robot’s orientation was kept constant throughout the image capturing process. We have taken this album of 300 images and unfolded all of them into rectangular panoramic images, 180×48 in size. Figure 3 shows an example image from the album and an unfolded version of the same image.

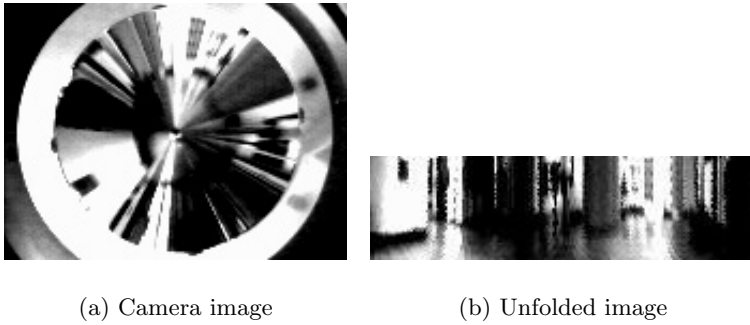


Fig. 3. The image unfolding process.

Figure 4 shows the homing performance of both the XALV and CGSM models on the image dataset using the image taken at position (210,180) as the snapshot image. This image is depicted in figure 3(b). Both models calculate home vectors for all images. These home vectors are depicted in figure 4. The figure also shows the success of homing for all grid positions. A pure white cell indicates successful homing. Homing is considered successful if a path along the home vectors exists from the position in question to the home position. We calculate a quantity called the *approach index* for each position,

$$\text{approach index} = \frac{\text{starting distance} - \text{closest distance}}{\text{starting distance}}$$

The quantity *starting distance* is the euclidean distance from the starting position to the home position. *Closest distance* is the smallest distance from the current position to the home position along the path from the starting position. The approach index for a successful starting position is 1. If, however, the path from the starting position leads away or perpendicular to the home position then the approach index will be 0. If the path approaches the home position but does not reach it then the approach index will be between 0 and 1. The purpose of the approach index is to judge the success of homing even if it is not completely successful. The level of whiteness shown in figure 4 is proportional to the approach index for the corresponding position.

Figure 4 has a dashed box surrounding the left half of the dataset. This box encloses the positions used to generate the results presented in [10]. It had been our intention to replicate these results precisely, however we have not been able to reconstruct all of the necessary parameters. This being said, our re-implementation of the XALV model seems to achieve better results. The original result for positions within the dashed box was that successful homing could be achieved from 98 of the 150 grid positions. Our re-implementation achieves successful homing from 136 of these 150 positions. This improvement is eclipsed by CGSM which achieves perfect homing (150 out of 150) within the dashed box.

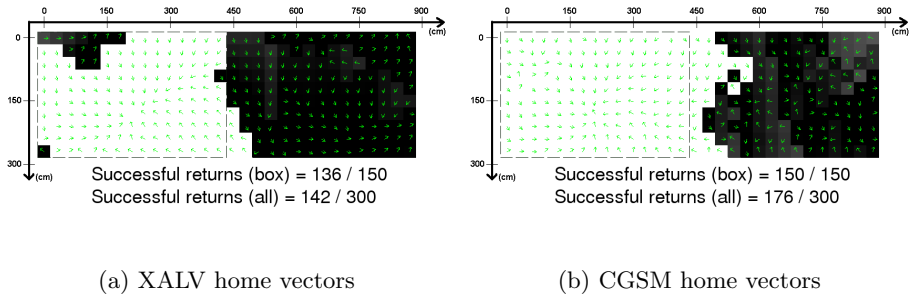


Fig. 4. Homing maps of the XALV and CGSM models for goal position (210,180). Home vectors are shown for images taken at corresponding positions. The whiteness of each position is proportional to that position’s approach index. Positions within the dashed box area were used for the results presented in [10].

Considering the whole of the dataset, our re-implementation of XALV achieves 142 / 300, while CGSM achieves 176 / 300. Both methods do relatively poorly in that part of the grid furthest from the home position.

The success of homing to one particular goal position does not tell the whole story. It may be that one method achieves good results for certain goal positions but not for others. We have generated 300 homing maps for the 300 image dataset using each image of the dataset, in turn, as the home position. For each homing map we count the number of successful homing positions and divide by 300. This is the *return ratio*. Figure 5 shows the return ratio for all images as the level of whiteness of each cell. Also shown is the approach index ratio which is the sum of all approach index values divided by 300.

The average and maximum return ratios are higher for CGSM than for XALV. The average and maximum approach index ratios are also higher. The standard deviation of both these ratios is lower for CGSM. For 268 of the 300 homing maps the return ratio is strictly higher for CGSM than for XALV. The approach index ratio is also strictly higher for 269 of 300 homing maps.

The area that surrounds a goal position from which successful homing can be achieved is known as the *catchment area*. We can translate the return ratio into squared metres to give an idea of the size of the catchment area for these two models. The average catchment area for CGSM is 14.6 m^2 (max. 25.3 m^2) while the average area for XALV is 6.7 m^2 (max. 21.5 m^2).

3.1 Discussion

Some caveats should first be mentioned in interpreting the results presented above. Firstly, all images were taken with constant orientation. This is an unrealistic scenario for actual robots homing in the real world. Even for robots with sophisticated odometry and compass systems, errors in orientation are inevitable.

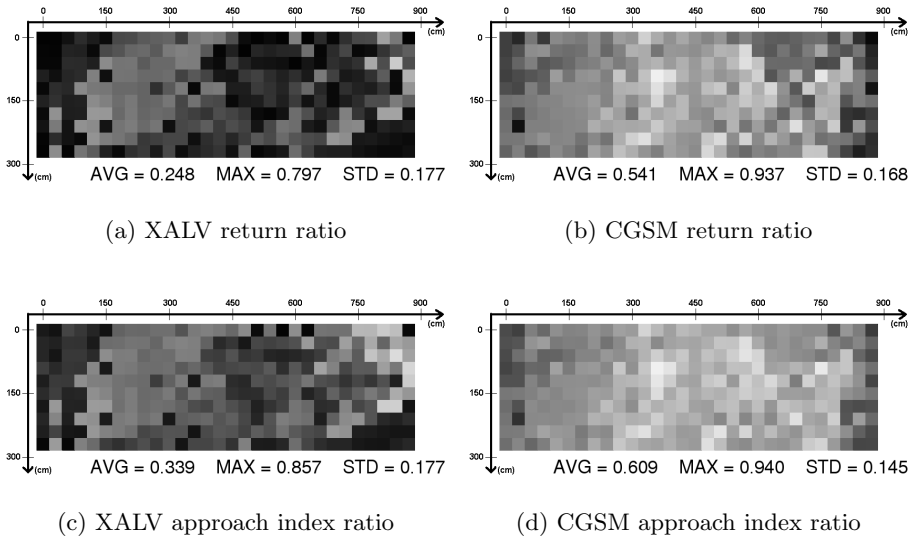


Fig. 5. Return ratio and approach index maps for XALV and CGSM.

Secondly, it was reported in [10] that changing lighting conditions seriously impacted on the homing performance of XALV. Thus, care must be taken in citing the catchment areas above as if they referred to actual environment-independent measures. Yet, it may also be true that the catchment area for this environment is *larger* than quoted above. If a set of images covering a larger area was employed it might be found that the catchment area would extend out further. Indeed in figure 4(b) the catchment area is prematurely bounded in places by the edges of the grid.

Our results show superior performance for CGSM over XALV. Both of these models can be considered variants of the snapshot model. We take the superior performance of CGSM to mean that snapshot-based homing models can do better by utilizing information that is implicitly encoded in both dimensions of an image seen by a homing agent. Further, we have shown that the success of a biologically plausible model tested only in simulation [8] can be transferred into a more realistic test environment where homing is based upon real-world images.

4 Conclusions

We have presented a variant of the snapshot model which is capable of successful homing using real-world test images. It exhibits improved performance over a model which did not make use of both input image dimensions. The structure and nature of our processing matrix was inspired by real biological systems. Areas for future work include free-roving robotic experiments, analysis of error

conditions, and the use of alternate features and feature movement detection methods.

Acknowledgments. Thanks to Ralf Möller and Thorsten Roggendorf for making their images available, assisting with the replication of their results, and for helpful comments. This work has been partially supported by NSERC Postgraduate Scholarship B - 232621 - 2000.

References

1. B. A. Cartwright and T.S. Collett. Landmark learning in bees. *Journal of Comparative Physiology*, 151:521–543, 1983.
2. T.S. Collett and J. Baron. Biological compasses and the coordinate frame of landmark memories in honeybees. *Nature*, 368:137–140, 1994.
3. Matthias O. Franz, Bernhard Schölkopf, Hanspeter A. Mallot, and Heinrich H. Bülthoff. Where did i take that snapshot? scene-based homing by image matching. *Biological Cybernetics*, 79:191–202, 1998.
4. J. Hong, X. Tan, B. Pinette, R. Weiss, and E.M. Riseman. Image-based homing. In *Proceedings of the 1991 IEEE International Conference on Robotics and Automation, Sacramento, CA*, pages 620–625, 1991.
5. D. H. Hubel and T.N. Wiesel. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *Journal of Physiology*, 160:106–154, 1962.
6. S.P.D. Judd and T.S. Collett. Multiple stored views and landmark guidance in ants. *Nature*, 392:710–714, 1998.
7. D. Lambrinos, R. Möller, T. Labhart, R. Pfeifer, and R. Wehner. A mobile robot employing insect strategies for navigation. *Robotics and Autonomous Systems, Special Issue: Biomimetic Robots*, 1999.
8. R. Möller, Marinus Marus, and D. Lambrinos. A neural model of landmark navigation in insects. *Neurocomputing*, 26-27:801–808, 1999.
9. Ralf Möller. Insect visual homing strategies in a robot with analog processing. *Biological Cybernetics*, 83(3):231–243, 2000.
10. Ralf Möller, Dimitrios Lambrinos, Thorsten Roggendorf, Rolf Pfeifer, and Rüdiger Wehner. Insect strategies of visual homing in mobile robots. In B. Webb and T. Consi, editors, *Biorobotics – Methods and Applications*. AAAI Press / MIT Press, 2001.
11. Thomas Röfer. Controlling a wheelchair with image-based homing. In *Proceedings of AISB Workshop on Spatial Reasoning in Mobile Robots and Animals, Manchester, UK*, 1997.
12. Thorsten Roggendorf. Visuelle Landmarkennavigation in einer natürlichen, komplexen Umgebung. Master's thesis, Universität Bielefeld, 2000.
13. E. Trucco and A. Verri. *Introductory Techniques for 3-D Computer Vision*. Prentice Hall, 1998.
14. Keven Weber, Svetha Venkatesh, and Mandyam Srinivasan. Insect-inspired robotic homing. *Adaptive Behavior*, 7:65–97, 1999.
15. J. Zeil, A. Kelber, and R. Voss. Structure and function of learning flights in bees and wasps. *Journal of Experimental Biology*, 199:245–252, 1996.