

1 Journal of Integrative Neuroscience, Vol. 4, No. 2 (2005) 1–18
 © Imperial College Press



3 **Research Report**

5 **EVENT BASED SELF-SUPERVISED TEMPORAL
 INTEGRATION FOR MULTIMODAL SENSOR DATA**

7 EMILIA I. BARAKOVA*

*RIKEN Brain Science Institute,
 2-1 Horosawa, Wako-shi, Saitama, 351-0198, Japan
 emilia@brain.riken.jp*

11 TINO LOURENS

*Honda Research Institute Japan Co., Ltd.,
 8-1 Honcho, Wako-shi, Saitama, 351-0114, Japan
 tino@jp.honda-ri.com*

15 Received 20 October 2004

Revised 18 April 2005

17 A method for synergistic integration of multimodal sensor data is proposed in this paper.
 18 This method is based on two aspects of the integration process: (1) achieving synergistic
 19 integration of two or more sensory modalities, and (2) fusing the various information
 20 streams at particular moments during processing. Inspired by psychophysical experiments,
 21 we propose a self-supervised learning method for achieving synergy with combined repre-
 22 sentations. Evidence from temporal registration and binding experiments indicates that dif-
 23 ferent cues are processed individually at specific time intervals. Therefore, an event-based
 24 temporal co-occurrence principle is proposed for the integration process. This integration
 25 method was applied to a mobile robot exploring unfamiliar environments. Simulations
 26 showed that integration enhanced route recognition with many perceptual similarities;
 27 moreover, they indicate that a perceptual hierarchy of knowledge about instant movement
 contributes significantly to short-term navigation, but that visual perceptions have bigger
 impact over longer intervals.

29 *Keywords:* Multimodal integration; robotics; navigation; proprioception.

31 **1. Introduction**

32 The world around us supplies huge amounts of information continuously from which
 33 living organisms extract the knowledge and awareness we need for survival. A fun-
 damental cognitive feature that makes this possible is the brain's ability to integrate
 all the various sensory inputs into a coherent representation of its environment. By

*Corresponding author.

2 *Barakova & Lourens*

1 analogy, robots are designed to continuously record large amounts of data using
various sensors, but their effectiveness suffers from a major design flaw. They lack
3 reference of how the information from the different sensory streams can be integrated
into consistent representations.

5 Multimodal integration for navigation has been studied from many differ-
ent perspectives. For instance, some studies have taken a purely computational
7 approach for their multimodal systems [19, 30]. Other studies took inspiration from
observable animal or human behaviors [3, 39, 24, 28, 41]. A third approach mod-
9 els neural mechanisms in areas of the brain associated with multimodal integration
[1, 2, 3, 11, 30, 31, 35, 38].

11 Our novel solution for this multimodal sensory integration problem for robot
navigation task was inspired by psychophysical experiments connecting visual and
13 idiothetic information. While navigating around an area, idiothetic information is
internally generated as the body moves through space [9, 25]. This information can
15 be derived from proprioceptive sensory streams about own movements and motor
efferent signals. Vestibular information, which follows the change in linear or rota-
17 tional movement velocity, is another source of idiothetic information. We integrated
view-based and velocity sensory information in our robotic implementations.

19 In general, two types of problems have to be solved for multimodal integration:
how and when should sensory cues be fused and how to obtain a synergistic mul-
21 timodal integration. Synergistic integration should produce more information from
the integrated representation than is evident from information generated in the sep-
23 arate modalities.

The first set of problems concerns the representational and technical aspects of
25 the actual fusion process. The method we proposed represents continuous sensory
information dynamically, by encoding the temporal history of sensor readings. This
27 encoding is a simple model of short-term memory. The method assumes that different
percepts unify in the brain, as suggested from temporal registration and binding
29 experiments showing that the information from one type of sensors is processed
separately on a certain time interval [10, 43]. To incorporate these observations into
31 a computational principle, we separate the processing of individual data streams
using a self-organizing principle until substantially different sensory information (or
33 a distinctive event) is perceived. Integration of the two sensory modalities takes place
only when the timing of a distinctive event encountered by both sensory streams
35 coincides. Since that distinctive event may occur at any moment in time, event-
based integration must divide the sensory streams into non-equal time intervals.
37 Event-based integration is a distinctive feature of the proposed integration method.

Synergy in a multimodal integration approach is difficult to quantitatively eval-
39 uate. Conflict estimation is a systematic approach that provides guidelines for
integration [39, 26, 40]. Conflict studies [39, 26, 40] investigate causes of spatial dis-
41 crepancies between the shifted spatial layout obtained through vision and the correct
spatial layout provided by other sensory modalities, such as proprioception. This
43 discrepancy (conflict) estimation approach has been used to judge spatial direction

1 perception as measured through target-pointing responses. It can be adapted to
2 evaluate how dynamic visual and non-visual information is integrated over time to
3 determine how distance traveled is perceived while moving.

4 One approach to integration is to add priority values (weights) to both (visual
5 and proprioceptive) information streams. Priorities are commonly assigned to visual
6 information [39, 26, 40], but the selection criteria are rarely disclosed and there are
7 differing opinions on how modality weights should be determined. One suggestion
8 is to base weights according to the precision of the information in each modality.
9 How the priority values are chosen, however, has not been shown by experimental-
10 ists. There are different ideas about how to determine the weights given to each
11 modality. According to one idea, the weights are determined by the precision of
12 the information in each modality [26, 42]; another assigns weights according to the
13 amount of attention directed to each modality [39, 20, 21, 40]. The thinking behind
14 these ideas stems from the concepts underlying statistical optimization models which
15 assume that sensory information from multiple sources should be weighted accord-
16 ing to the estimated reliability of each cue. Unlike the discussed issues surrounding
17 weight assignments for vision and proprioception cues in conflict studies [39, 26, 40],
18 we coupled the dynamic view-like and self-motion information using a self-teaching
19 principle, conforming with the application domain of concurrent mapping and nav-
20 igation [4, 11, 45].

21 To illustrate our integration approach, two data streams were recorded as an
22 autonomous robot explored an unfamiliar environment. They provided absolute and
23 relative information about the robot's movement with respect to the relation of robot
24 movement to the surrounding objects.

25 In Sec. 2, we present our integration hypothesis followed by an explanation of
26 the temporal synchronization principle as the framework for the application domain
27 in Sec. 3. The actual integration and the experimental testing is presented in Sec. 4
28 and the results are discussed in Sec. 5.

29 **2. Hypothesis**

30 Multisensory integration requires an anatomical convergence of unisensory inputs
31 onto a single neuron or ensembles of interconnected neurons [33], and some degree
32 of temporal alignment of the unisensory inputs [36]. Areas associated with multisen-
33 sory integration include the superior temporal polysensory area, lateral and ventral
34 intraparietal areas. A detailed explication of the brain mechanisms of multisensory
35 processing has been conducted in the carnivore superior colliculus [36] and substan-
36 tial progress has also been made at the neocortical level, most notably in monkeys
37 [5, 7, 13, 14, 16, 17] and recently in humans [8, 18].

38 Despite these advances, questions remain about the anatomical substrates of
39 multisensory convergence in primates. Questions also remain about the temporal
40 parameters of the converging sensory inputs. Temporal windows exist for the inte-
41 gration of neural responses to stimulus inputs from different modalities and for

1 perception of fused multisensory inputs (i.e., relating to the same object [36]). How-
ever, only a few studies of the timing of sensory inputs to the neocortex have been
3 reported, focusing mostly on the response latencies in the visual system [26, 32, 30].
Therefore, detailed modeling at this stage remains difficult.

5 Temporal registration experiments suggest that the brain does not bind infor-
mation entities from different modalities in real time; instead, it binds the results
7 of its own processing systems on certain time intervals. Our hypothesis provides a
constructive basis for an integration strategy. We assumed that the brain operates
9 as a self-organizing information system that processes the sensory flow to asymmet-
ric activation patterns at the separate sensory modality level. Since each modality
11 brings different levels of generalization and information about the external world,
information from any one modality can also serve as a “teacher” for other modalities.

13 Experimental support for this line of thinking comes from studies of the rat
head-direction system. Head-direction cells are located in many parts of the rat
15 brain, including the pre- and post-subiculum, the anterodorsal thalamus, and the
mammillary nuclei [27, 30, 34, 35]. Head-direction tuning may arise when angular
17 head-velocity signals are fused. These signals originate from vestibular neurons that
are tonically active when the head is still, but the firing rate increases when turning
19 is in one direction (i.e., center to left) and decreases when it turns the other way
(i.e., center to right).

21 Experimental evidence also implicates a fast-acting contribution of visual input
in the organization of head-direction circuits [6, 8]. That is, the preferred head direc-
23 tion of these cells can be controlled by a visually salient landmark. When the head
is rotated, the preferred direction of these cells is generally aligned with the angular
25 displacement of the landmark [37]. These, and other similar experiments [30, 46],
suggest a critical role of visual information in the calibration and development of
27 head-direction tuning. Calibration of head-direction cells by visual landmarks has
been shown in [44]. The possible role of visual information as a teaching signal that
29 supervises the development of an integrator network has also been studied [15]. They
concluded that selective amplification teaches the vestibular input how to predict
31 and replace missing input.

3. Temporal Synchronization in Spatial Navigation Setup

33 In robots, various sensors asynchronously provide information with different mean-
ings and sampling characteristics. Established ways of combining the information
35 from different information sources are missing. To combine multimodal information
sources, the following principles need to be considered:

- 37 • Data that are perceived (recorded) at the same time interval relate to the same
situation (event).
- 39 • Processing of different data streams is done in separate modalities, followed by
synchronization using a temporal principle.
- 41 • The temporal synchronization is event-based (not fixed-time interval based).

1 With these guiding principles, the following computational steps will be taken:
2 firstly, event-based time intervals will be defined. Secondly, information recorded
3 within these intervals will be represented so that fusion is possible. Thirdly, the
4 actual integration takes place.

5 The mapping task is solved by using the data recorded by the robot during its
6 exploration of an unknown environment. Figure 1 shows the experimental environ-
7 ment. Several exploration routes are plotted on the picture. Black points on the floor
8 indicate novel regions that have been clustered in different classes (events) according
9 to the sensory information as detailed below.

10 Imagine a group of sensors imitating a consummate description of the environ-
11 ment that can be created by biological systems is difficult. A plausible alternative
12 might be to focus on a couple of sensors that provide complementary informa-
13 tion about the environment. The relevance of an egocentric perspective for an
14 autonomous robot in spatial modeling of previously unknown environments was
15 elaborated in [4, 45]. That egocentric model combines two types of information:
16 absolute and relative with respect to the relation of robot movement to the sur-
17 rounding objects. The absolute perspective records sensory information independent
18 of robot movements using laser range finders. Relative information reflects the robot
19 perceptions of its motion. For instance, if the flooring is different, the angular vel-
20 ocity readings may differ when the robot takes the very same trajectory. The relative
21 sensory stream is recorded by a build-in gyroscope.

22 The “views” that the robot perceives with a laser range finder source the absolute
23 information. Visual information has an absolute character. Since this method takes

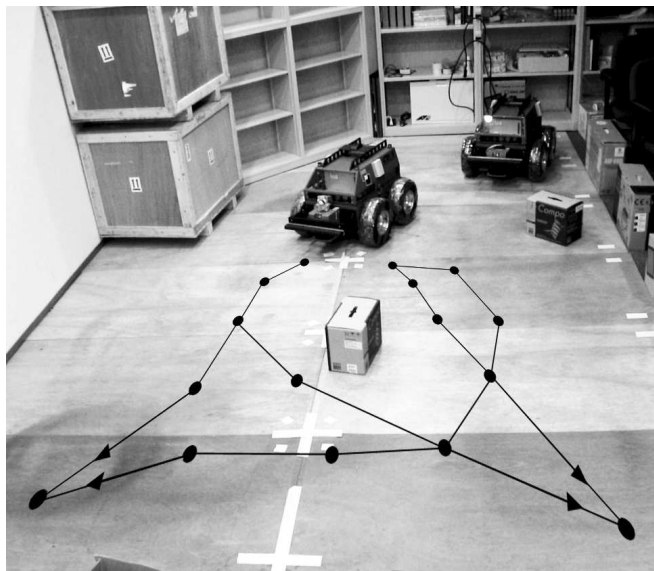


Fig. 1. Experimental environment in the picture several recorded trajectories that reach perceptually similar places have been added.

6 *Barakova & Lourens*

1 ideas from the combination of visual and proprioceptive cues by living organisms,
 2 the recordings from the laser range finder are referred to as view-like information.
 3 Vision is not used because the application domain that the current study is target-
 4 ing is underwater navigation, an area where visual information is scarcely available.
 5 Apparently, the range sensor information is simpler to process and sufficient to repre-
 6 sent the idea of fusing absolute with relative egocentric information. The individual
 7 “view” of the robot is formed by recording 360 samples per 180 degrees. A snapshot
 8 of a polar representation of such a record is shown in Fig. 2.

9 A sequence of such snapshots, recorded during robot exploration, form a dynamic
 trajectory:

$$11 \quad x_i(t_i) = \frac{s_i(t_i) + \sum_{\tau=1}^{h_i-1} f_i(\tau)x_i(t_i - \tau)}{h_i} \quad (3.1)$$

12 where $s_i(t_i)$ is the readout of sensor (element) i , $t_i > 0$ is an integer time stamp
 13 for sensor (element) i , $x_i(t_i)$ specifies the current sensory representation, that keeps
 14 a history of $h_i = \min(H_i, t_i)$ elements from the exponentially decaying forgetting
 15 curve f . A priori known constant H_i denotes the maximum history length. Items in
 16 the dynamic trajectory decay in time, corresponding to the decay theory of forgetting
 17 in short term memory. In our experiments a decay kernel [Eq. (3.2)] was used. This is
 because, of all the previously seen patterns, the last is the most vivid and influences

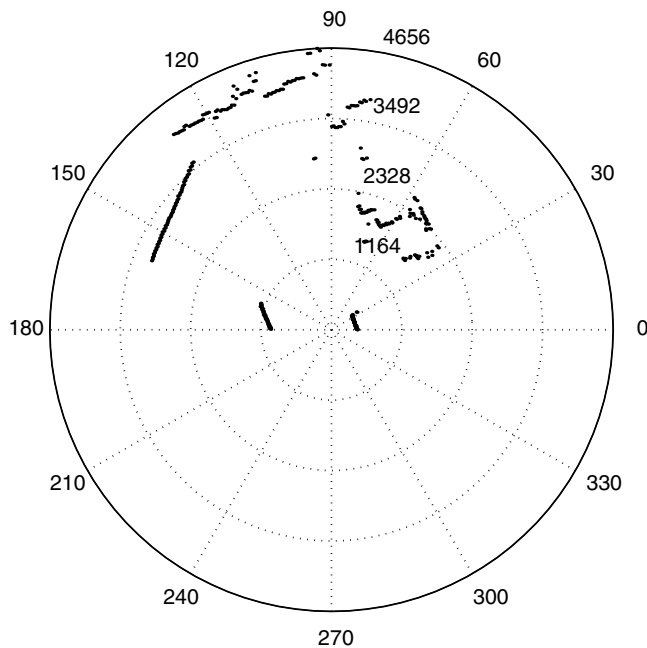


Fig. 2. Sample recording from the laser range finder. The range finder “view” is a composite recording of 360 samples per 180 degrees that forms a snapshot. Snapshots are recorded at frequency of 4.7 Hz. The distances are presented in millimeters.

1 most strongly the current perception, as the influence of the older patterns fade
away.

$$3 \quad f_i(t) = \exp(-\alpha_i t) \quad (3.2)$$

where constant $\alpha_i \in [0, 1]$ determines the decay profile f_i for sensor (element) i . Each
5 unit in this short-term memory model samples a symbol in a specific time interval.
Such a dynamic sequence encodes the first sensory stream, used for the integration
7 (i.e., the laser range finder stream). It represents the absolute perspective of the
robot about its own motion (i.e., robot motion with respect to the surrounding
9 objects).

The velocity measured by the gyroscope is used to represent the relative cue that
11 resembles proprioceptive information in animals. It reflects the robot's perception of
its own motion. Most informative are the angular velocity recordings from the robot,
13 since they reflect directional changes in its trajectory. The temporal synchronization
of these two information streams is performed as follows: a neural gas algorithm
15 [23] is applied to the view-based sensory stream to determine the intervals when a
novel "event" occurs. Learning dynamics is guided by a combination of competitive
17 Hebbian learning and vector quantization algorithm. A set of n reference vectors \vec{w}_i ,
 $i \in \{1, 2, \dots, n\}$ have strengths, depending on their neighborhood ranking. When an
19 input vector $\vec{x} = \{x_1(t_1), x_2(t_2), \dots, x_n(t_n)\}$ is presented, a neighborhood ranking
of the reference vectors takes place ($\vec{w}_{i0}, \vec{w}_{i1}, \dots, \vec{w}_{in}$) with \vec{w}_{i0} being the closest to
21 \vec{x} , \vec{w}_{i1} being the second closest to \vec{x} , and \vec{w}_{ik} , $k \in \{1, 2, \dots, n-1\}$ is the reference
vector for which there are k vectors \vec{w}_j with $\|\vec{x} - \vec{w}_j\| < \|\vec{x} - \vec{w}_{ik}\|$. The rank index
23 associated with \vec{w}_i is denoted by $k_i(\vec{x}, \vec{w})$. Using a Hebbian-like rule, the adaptation
step for adjusting is given by:

$$25 \quad \Delta \vec{w}_i = \varepsilon(t) h_\lambda(k_i(\vec{x}, \vec{w})) (\vec{x} - \vec{w}_i) \quad (3.3)$$

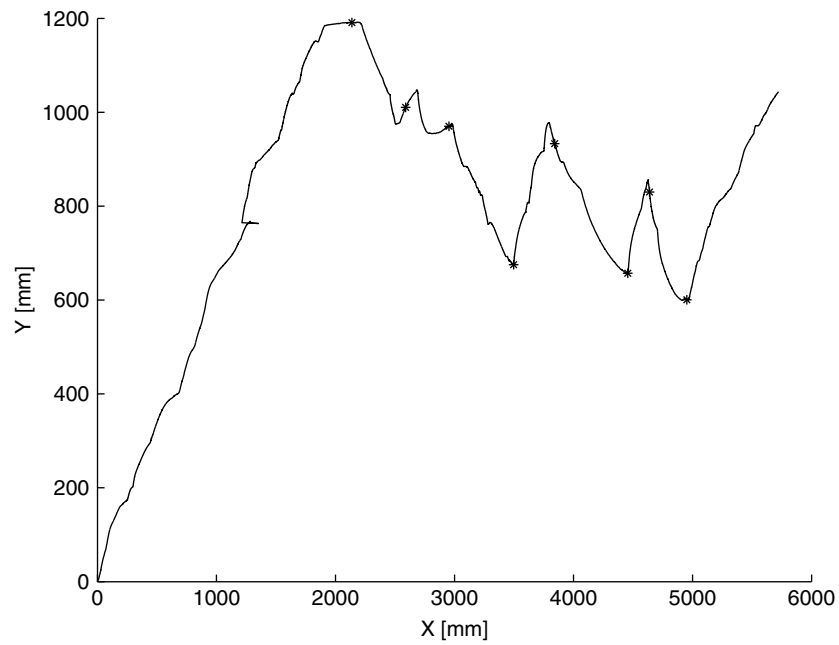
where the step size $\varepsilon \in [0, 1]$ is the learning rate, and $h_\lambda(k_i(\vec{x}, \vec{w})) \in [0, 1]$ accounts
27 for the topological arrangement of the \vec{w}_i in the input space.

$$h_\lambda = \exp\left(-\frac{k}{\lambda}\right) \quad (3.4)$$

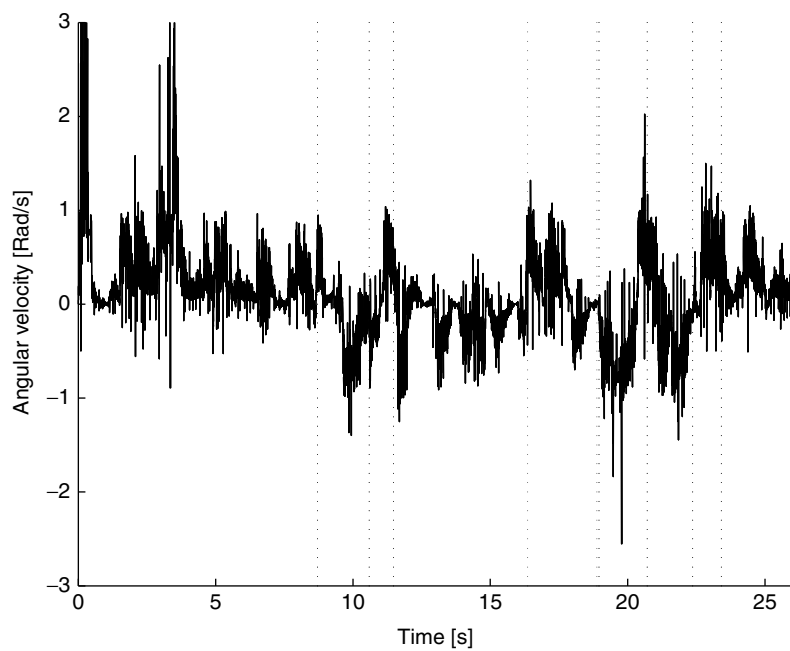
29 i.e., the neighborhood relies on the rank in the ordered sequence of distances, and
the weights are learned according to (3.3), while decreasing λ . For the simulation,
31 the $\varepsilon(t)$ and $\lambda(t)$ are calculated as follows:

$$g(t) = g_0 \left(\frac{N}{g_0}\right)^{\frac{t}{T}} \quad (3.5)$$

33 where $g \in \{\lambda, \varepsilon\}$, $N = 0.01$, $\varepsilon_0 = 0.5$, $\lambda_0 = n/2$, where n is the number of neurons,
and T is the number of the training patterns. The simulations were made with
35 $n = 20$ neurons, sufficient to encode the different patterns from the experimental
environment. In order to extend the method to any environment, an incremental
37 version of the algorithm known as a growing neural gas algorithm [12] can be used.

8 *Barakova & Lourens*

(a)



(b)

Fig. 3. Temporal synchronization principle. (a) Experimental trajectory is segmented by the clustering algorithm. Points on the curve denote the events that are determined as new. (b) The corresponding angular velocity curve, segmented according to the time the robot has spent in the same event cluster.

1 This algorithm starts from two neurons and increases the neurons number, until all
the different patterns are encoded to separate classes. A Euclidean distance measure
3 decides how many classes to form.

4 Synchronously, an event-based segmentation is performed on the second (veloc-
5 ity) sensory stream. Note that the dynamic events are formed on-line with minimal
processing of incoming data.

7 Figure 3 illustrates the temporal synchronization of the two information streams.
Figure 3a represents a trajectory that the robot took during its exploration of an
9 environment. Any qualitatively different “view” that the robot observes is defined
as new segment in the environment and it is denoted by * in Fig. 3. The duration of
11 the trajectory in these segments determines the division of the other sensory data
stream as illustrated in Fig. 3b.

13 The results of a temporal synchronization by a recorded short experimental tra-
jectory are shown in Fig. 4. Figure 4a shows view-based segmentation of a two-
15 dimensional space, obtained during free exploration of the robot. The points where a
qualitatively new view occurs, as determined by the described algorithm, are shown.
17 The robot trajectory is not shown on the plot. Corresponding velocity trajectory is
reconstructed by angular and linear velocity recordings (Fig. 4b). Figure 4c shows
19 the velocity curve after temporal synchronization with the view clusters.

21 The synchronization process is as follows: after clusters of the view based infor-
mation stream are found, the velocity data are segmented on the same time inter-
23 vals, considering the different sampling frequencies of both sensory streams. Based
on this segmentation, a clustering to unified trajectory elements (or motion prim-
25 itives) of the velocity curve is made, using the algorithm described in Eqs. (3.3)–
(3.5). The input vectors by this clustering are the segments from the velocity curve.
27 Figures 4d–f depict the resulting motion primitives corresponding to the shown seg-
ments of the velocity curve.

4. Integration Results

29 In the learning phase, the sequences of views were clustered by the self-organizing
algorithm, forming dynamic events throughout the robot’s continuous exploration
31 of the unknown environment. Those moments when a qualitatively new view is
recorded by the gyroscope become the temporal dividers for the proprioceptive
33 dynamic sensory stream. Similar steps in the testing phase were performed with the
sensory data gathered from further exploration. Specific trajectories were recorded
35 for the testing process where classification, not clustering, is done based on the
clustering categories already attained from the exploration data. The database of
37 distinctive events is built gradually during on-line implementation. Figure 5 summa-
rizes the computational flow by off-line event based integration. A supervised neural
39 network is used for teaching between the modalities, that are denoted as multimodal
integration blocks in the figure.

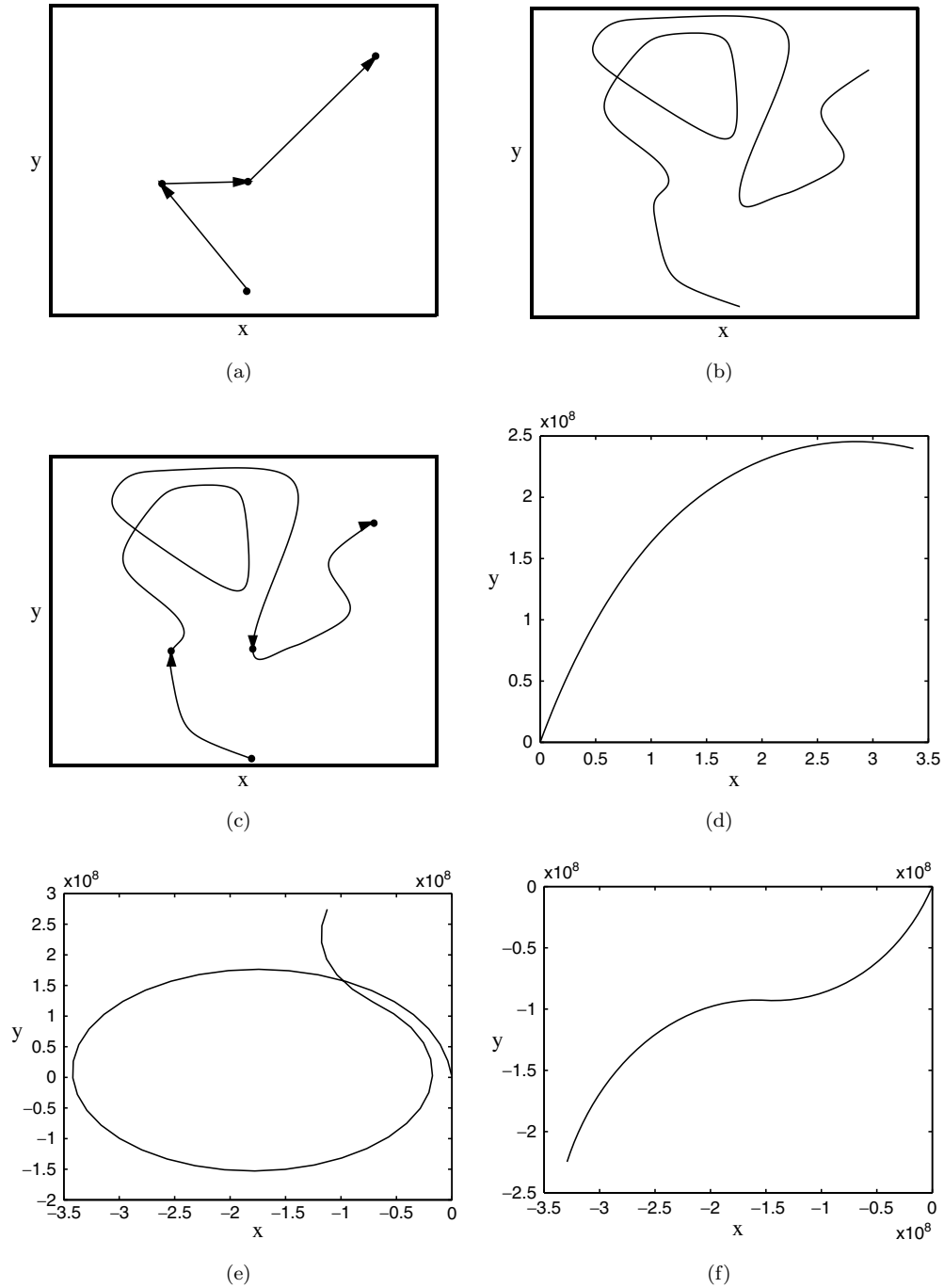


Fig. 4. The result of temporal synchronization over a short trajectory recorded by the robot. (a) View-based segmentation of a two dimensional space. (b) The corresponding velocity trajectory (c) Temporally synchronized velocity curve. (d-f) Classes, corresponding to the segments of the velocity curve.

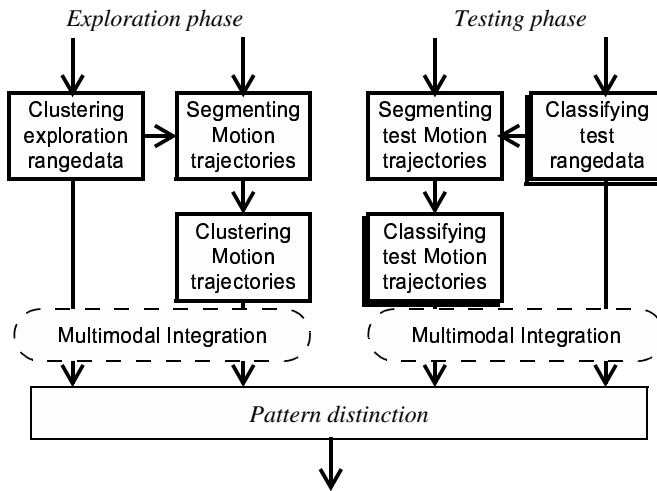


Fig. 5. Information flow of event-based integration, consisting of exploration and testing phases that reflect the experimental process.

1 Several groups of experiments were performed according to this computational
 2 scheme. In the first group of experiments, perceptually similar places with respect to
 3 the view-based sensory stream were training inputs and motion primitives perceived
 4 by the auto-motion (velocity) stream functioned as teacher. Figure 6 depicts the
 5 dynamic trajectories, obtained from the view-based information stream.

6 All trajectories reach one of two corners in the experimental environment, which
 7 look the same if “seen” by the laser range sensors (Fig. 6a). Figures 6b–f depict the
 8 encoding by dynamic trajectory formation. The plots on the left show the dynamic
 9 trajectory and the plots on the right show the classes of distinctive views encoun-
 10 tered while making these trajectories. Considering only the temporal history of
 11 experienced views, the dynamic trajectory method can distinguish most of the tra-
 12 jectories, although two, shown in Figs. 6e–f, appear similar despite the obviously
 13 different view history. Using the self-motion primitives (that can be obtained from
 14 the velocity sensory cue) as a teaching signal, all dynamic trajectories are disam-
 15 biguated. Experimental analysis shows that using self motion as a teaching signal
 16 helps perceptual aliasing, but sometimes distinguishes trajectories that are very
 17 similar.

18 The view-based signal was used as teacher and the self motion signal as an input
 19 in the second group of experiments. New motion trajectories, specifically selected for
 20 their similarity with respect to classification to motion primitives, were trained. The
 21 obtained results coincided better with decisions that a human observer would make.

22 The later experiments were made with trajectories that were not controlled with
 23 respect to similarities. Correspondingly, short, long, and arbitrary length sequences
 24 were recorded. The environment contains 6 big objects that divide it on routes
 25 that can be distinguished by a human observer. All possible trajectories between

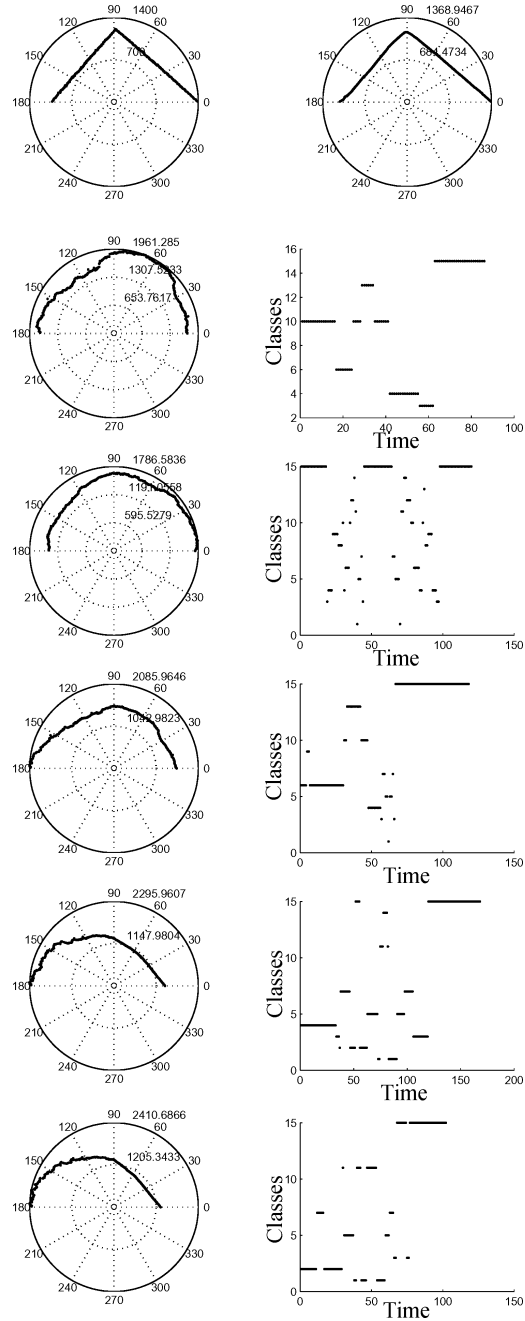
12 *Barakova & Lourens*

Fig. 6. Environmental similarities and the corresponding dynamic trajectories. (a) The left plot shows the laser range finder recording of a corner in the experimental room. The plot on the right shows the network output. (b–f) 5 dynamic trajectories that finish at that corner (i.e., have the same final sensory reading). Alternative routes to the corner are shown as clusters in the right plots.

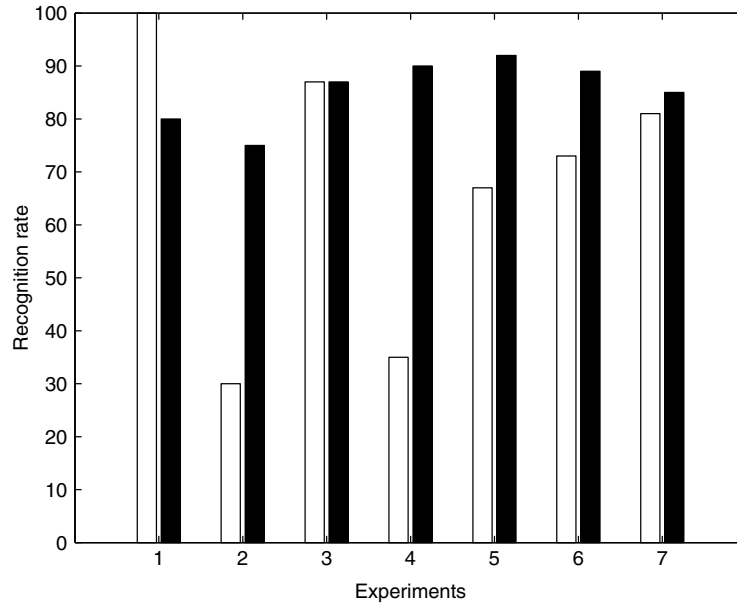


Fig. 7. Recognition rate with view based information used as a teacher (black bars) and self-motion information used as a teacher (white bars). The experiments are made as follows: 1-perceptually similar trajectories with respect to view based sensors; 2-perceptually similar trajectories as perceived by self-motion sensors; 3-short trajectories; 4-long trajectories; 5-7 non selected trajectories.

1 the same objects or object traversed in the same order belong to the same route.
 For the learning phase the robot has passed all the possible routes between the
 3 environmental objects at least once. The test sequences were recorded by a robot
 traversing trajectories within the possible routes. Every trajectory record consists
 5 of view sequence of data and velocity sequences for reconstructing the self-motion
 trajectories. The recognition of a trajectory was evaluated manually.

7 Experimental results where self-motion signals were used to teach are shown in
 Fig. 7 as white bars. Results with the view signal as teacher are plotted with black
 9 bars in the same figure. The recognition rate was evaluated as the percentage of the
 trajectories from the given group that was classified in the right class.

11 Classification over the same training sets was also made with the weighting
 integration method, implemented after [39]. The bigger weights are assigned to view
 13 information; however, better results with current method, as shown in Fig. 8, are
 not achieved even after thorough experimental testing of both algorithms. They are
 15 validated for this particular case only.

5. Discussion

17 The proposed multimodal integration method provides an alternative way to com-
 19 bine information from different sensory modalities. It addresses two groups of
 problems — firstly, how and when the combination between the two sensory streams

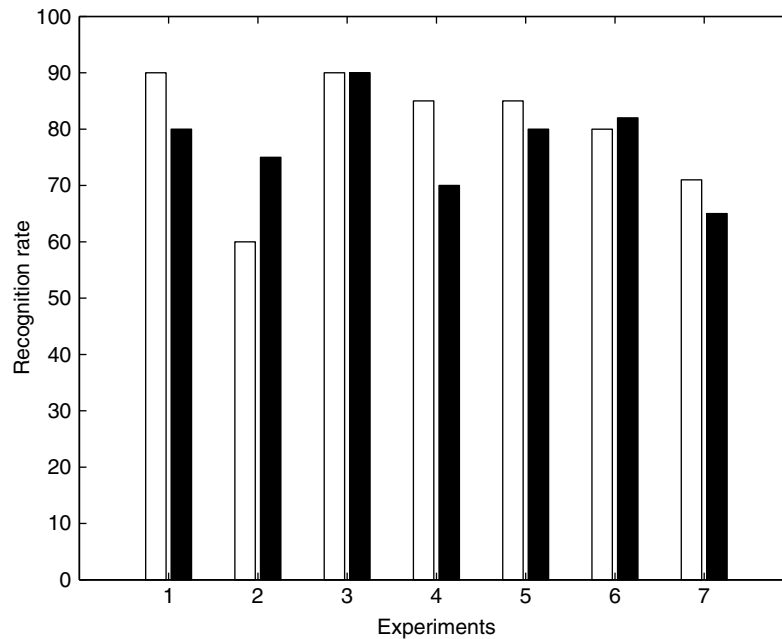


Fig. 8. Recognition rate by weighted integration (black bars) and integration with a teacher (white bars). The experiments are made as follows: 1-perceptually similar trajectories with respect to view based sensors; 2-perceptually similar trajectories as perceived by self-motion sensors; 3-short trajectories; 4-long trajectories; 5-7 non selected trajectories.

1 occurs, and secondly, how synergy in the integration process is obtained. The estab-
 2 lished synchronization framework is an efficient way to bring the multimodal infor-
 3 mation together at event-based temporal intervals. This approach prevents error
 4 accumulation since the synchronization takes place over relatively short temporal
 5 intervals (i.e., the accumulative errors are reset after every interval). Another advan-
 6 tage of this model is that both the consequent perceptions as well as the transitions
 7 between them are dynamically encoded.

8 Information from both sensory streams is used directly, without feature extrac-
 9 tion. The basis for this synergistic integration is tested for data sets recorded on
 10 the following principles: perceptual similarity of the trajectories with respect to
 11 view-based information; perceptual similarity of the trajectories as perceived by
 12 self-motion sensors; different length of the trajectories. The experimental results
 13 show that in most cases using a view-based sensory stream as a teacher is advan-
 14 tageous. Trajectories, recorded on a visual similarity principle can be successfully
 15 disambiguated if the self-motion (velocity) signal is used as a teacher, but this dis-
 16 ambiguation is crude, occasionally not differentiating between similar trajectories.
 17 Disambiguation ensured by the view-based teacher never reaches 100 percent recog-
 18 nition of experienced trajectories, but gives results that most closely mimic human
 19 observer decision behavior.

1 This integration method is compared to a weighted integration method that
2 uses visual and proprioceptive sensory cues (see [39]). This is possible since two sen-
3 sory streams are used: view-based information resembles animal vision, and velocity
4 information is equivalent to proprioception by living organisms in this experimen-
5 tal setting. In most of the testing set groups, the integration process gives better
6 results by the method, as proposed in this paper, except when there are similarities
7 in self-motion sensory stream or when the data sets are absent.

8 From the third and fourth groups' test sets we can conclude that the proposed
9 integration method explains perceptual hierarchy in the following way: knowledge
10 about instant movement contributes significantly to short-term navigation, while
11 visual perceptions have bigger impacts over longer terms.

12 While temporal synchronization can be implemented as an on-line learning pro-
13 cess, the multimodal integration requires off-line processing. This is one major draw-
14 back of the model.

15 We base our multimodal integration on studies that rely on visual and pro-
16 prioceptive data streams. Such streams have different dimensionality. However, the
17 experiments show integration of two unidimensional data streams. Using visual infor-
18 mation is a straightforward extension to the developed integration method, since the
19 two integration streams are processed separately. The used algorithm can work with
20 two-dimensional data as well. The speed of processing, however, will be a problem
21 for a real robotic implementation. To deal with that, we have the following idea.
22 The range sensor readings give a sparse depth representation. The same could be
23 obtained by using unimodal vision features. For instance the most salient features
24 of corners/crossing or edges could be used to construct a depth map [22].

25 Further, we plan to model the integration process on the level of cortical struc-
26 tures, including temporal synchronization with feedback, since recent work suggests
27 that feedforward and feedback projections contribute to the convergence of the inte-
28 grated representation. We will use visual and proprioceptive cues.

29 Acknowledgments

30 I (E.B.) would like to acknowledge the multiple discussions with U.R. Zimmer
31 (Australian National University, Canberra) for the multiple discussions. We have
32 done some initial experimental work together, see [4].

33 References

- 34 [1] Andersen RA, Multimodal integration for the representation of space in the posterior
35 parietal cortex, *Philos T Rly Soc B* **352**(1360):1421–1428, 1997.
- 36 [2] Babeau V, Gaussier P, Joulain C, Revel A, Banquet J, Merging visual place recognition
37 and path integration for cognitive map learning, *SAB*, 101–110, 2000.
- 38 [3] Banquet J, Gaussier P, Quoy M, Revel A, From reflex to planning: Multimodal, versa-
39 tile, complex systems in biorobotics, *Behav Brain Sci* **24**(6):1051–1053, 2001.

16 *Barakova & Lourens*

- 1 [4] Barakova EI, Zimmer UR, Global spatial modeling based on dynamics identification
 3 according to discriminated static sensations, in *Proceedings of Underwater Technology*,
 IEEE, Tokyo, 2000.
- 5 [5] Benevento LA, Fallon J, Davis BJ, Rezak M, Auditory-visual interaction in single cells
 in the cortex of the superior temporal sulcus and the orbital frontal cortex of the
 macaque monkey, *Exp Neurol* **57**:849–872, 1977.
- 7 [6] Blair H, Sharp P, Anticipatory head direction signals in anterior thalamus: evidence
 for a thalamocortical circuit that integrates angular head motion to compute head
 9 direction, *J Neurosci* **15**:6260–6270, 1995.
- [7] Bruce C, Desimone R, Gross CG, Visual properties of neurons in a polysensory area
 11 in superior temporal sulcus of the macaque, *J Neurophysiol* **46**:369–384, 1981.
- [8] Canon LK, Directed attention and maladaptive “adaptation” to displacement of the
 13 visual field, *J Exp Psychol* **88**:403–408, 1971.
- [9] Chance SS, Gaunet F, Beall AC, Loomis JC, Locomotion mode affects the updating
 15 of objects encountered during travel: the contribution of vestibular and proprioceptive
 inputs to path integration, *Presence-Teleop Virt* **7**:168–178, 1998.
- 17 [10] Eagleman DM, Sejnowski TJ, Motion integration and postdiction in visual awareness,
Science **287**(5460):2036–2038, 2000.
- 19 [11] Elfes A, Sonar-based real-world mapping and navigation, *IEEE J Robot Autom*
RA-3(3):249–265, June 1987.
- 21 [12] Fritzke B, A growing neural gas network learns topologies, *Adv Neural Inf Process Syst*
 625–632, 1995.
- 23 [13] Graziano MS, Hu XT, Gross CG, Visuospatial properties of ventral premotor cortex,
J Neurophysiol **77**:2268–2292, 1997.
- 25 [14] Graziano MS, Yap GS, Gross CG, Coding of visual space by premotor neurons, *Science*
266:1054–1057, 1994.
- 27 [15] Hahnloser RHR, Emergence of neural integration in the head-direction system by visual
 supervision, *Neuroscience* **120**(3):877–891, 2003.
- 29 [16] Hikosaka K, Iwai E, Saito H, Tanaka K, Polysensory properties of neurons in the ante-
 31 rior bank of the caudal superior temporal sulcus of the macaque monkey, *J Neurophysiol*
60:1615–1637, 1988.
- [17] Hyvarinen J, Shelepin Y, Distribution of visual and somatic functions in the parietal
 33 associative area 7 of the monkey, *Brain Res* **169**:561–564, 1979.
- [18] Jousmaki V, Hari R, Parchment-skin illusion: sound-biased touch, *Curr Biol* **8**(6):R190,
 35 1998.
- [19] Kam M, Zhu X, Kalata P, Sensor fusion for mobile robot navigation, *P IEEE* **85**(1),
 37 1997.
- [20] Kelso JAS, Cook E, Olson ME, Epstein W, Allocation of attention and the locus of
 39 adaptation to displaced vision, *J Exp Psychol Human* **1**:237–245, 1975.
- [21] Loomis JM, Da Silva JA, Fujita N, Fukusima SS, Visual space perception and visually
 41 directed action, *J Exp Psychol Human* **18**:906–921, 1992.
- [22] Lourens T, Würtz RP, Extraction and matching of symbolic contour graphs, *Int J*
 43 *Pattern Recogn* **17**(7):1279–1302, 2003.
- [23] Martinetz TM, Berkovich SG, Schulten KJ, Neural gas network for vector quantization
 45 and its application to time-series prediction, *IEEE T Neural Networ* **4**:558–569, 1993.

- 1 [24] Menzel R, Geiger K, Chittka L, Joerges J, Kunze J, Ller U, The knowledge base of bee
navigation, *J Exp Biol* **199**(Pt 1):141–146, 1996.
- 3 [25] Mittelstaedt ML, Mittelstaedt H, Idiopathic navigation in humans: estimation of path
length, *Exp Brain Res* **139**:318–332, 2001.
- 5 [26] Pick HL, Warren DH, Hay JC, Sensory conflict in judgments of spatial direction, *Per-*
cept Psychophys **6**:203–205, 1969.
- 7 [27] Ranck JJ, Head-direction cells in the deep cell layers of the dorsal presubiculum in
freely-moving rats, *Soc Neurosci* **10**:599, 1984.
- 9 [28] Rieser JJ, Ashmead DH, Talor CR, Youngquist GA, Visual perception and the guidance
of locomotion without vision to previously seen targets, *Perception* **19**:675–689, 1990.
- 11 [29] Roumeliotis S, Bekey G, An extended kalman filter for frequent local and infrequent
global sensor data fusion, *Proc SPIE* **3209**, 1997.
- 13 [30] Save E, Poucet B, Involvement of the hippocampus and associative parietal cortex in
the use of proximal and distal landmarks for navigation, *Behav Brain Res* **109**(2):195–
15 206, 2000.
- 17 [31] Schlack A, Hoffmann KP, Bremmer F, Interaction of linear vestibular and visual stim-
ulation in the macaque ventral intraparietal area (vip), *J Neurosci* **21**:23–29, 2001.
- 19 [32] Schmolesky M, Wang Y, Hanes DP, Thompson KG, Leutgeb S, Schall JD,
Leventhal AD, Signal timing across the macaque visual system, *J Neurophysiol*
79:3272–3278, 1998.
- 21 [33] Schroeder CE, Mehta AD, Foxe JJ, Determinants and mechanisms of attentional mod-
ulation of neural processing, *Front Biosci* **6**:672–684, 2001.
- 23 [34] Seltzer B, Pandya DN, Afferent cortical connections and architectonics of the superior
temporal sulcus and surrounding cortex in the rhesus monkey, *Brain Res* **149**:1–24,
25 1978.
- 27 [35] Seth AK, McKinsty JL, Edelman GM, Krichmar JL, Visual binding through reen-
trant connectivity and dynamic synchronization in a brain-based device, *Cereb Cortex*
14:1185–1199, 2004.
- 29 [36] Stein BE, Meredith MA, *The Merging of the Senses*, MIT Press, Cambridge, MA, 1993.
- 31 [37] Taube J, Muller R, Ranck J, Head-direction cells recorded from the postsubiculum in
freely moving rats: Ii. the effects of environmental manipulations, *J Neurosci* **10**:436–
447, 1990.
- 33 [38] Tononi G, Sporns O, Edelman GM, Reentry and the problem of integrating multiple
cortical areas: Simulation of dynamic integration in the visual system, *Cereb Cortex*
35 **2**:310–335, 1992.
- 37 [39] Van Beers RJ, Sittig AC, Denier van der Gon JJ, Integration of proprioceptive
and visual position-information: an experimentally supported model, *J Neurophysiol*
81:1355–1364, 1999.
- 39 [40] Warren DH, Schmitt TL, On the plasticity of visual-proprioceptive bias effects, *J Exp*
Psychol Human **4**:302–310, 1978.
- 41 [41] Wehner R, Michel B, Antonsen P, Visual navigation in insects: Coupling of egocentric
and geocentric information, *J Exp Biol* **199**(1):129–140, 1996.
- 43 [42] Welch RB, Widawski MH, Harrington J, Warren DH, An examination of the relation-
ship between visual capture and prism adaptation, *Percept Psychophys* **25**:126–132,
45 1979.

18 *Barakova & Lourens*

- 1 [43] Zeki S, Functional specialization in the visual cortex of the rhesus monkey, *Nature*
274:423–428, 1978.
- 3 [44] Zhang K, Representation of spatial orientation by the intrinsic dynamics of the head-
direction cell ensemble: a theory, *J Neurosci* **16**:2112–2126, 1996.
- 5 [45] Zimmer UR, Robust world-modelling and navigation in a real world, *Neurocomputing*
13:247–260, 1996.
- 7 [46] Zugaro M, Berthoz A, Wiener S, Background, but not foreground, spatial cues are
taken as references for head-direction responses by rat anterodorsal thalamic neurons,
9 *J Neurosci* **21 RC145**:1–5, 2001.