

Adaptive computational models of fast learning of motion direction discrimination

V. Sundareswaran, Lucia M. Vaina

Intelligent Systems Laboratory, College of Engineering, Boston University, Boston, MA 02215, USA

Received: 27 June 1995/Accepted: 21 October 1995

Abstract. In a previous study, we found that subjects' performance in a task of direction discrimination in stochastic motion stimuli shows fast improvement in the absence of feedback and the learned ability is retained over a period of time. We model this learning using two unsupervised approaches: a clustering model that learns to *accommodate* the motion noise, and an averaging model that learns to *ignore* the noise. Extensive simulations with the models show performance similar to psychophysical results.

1 Introduction

Perceptual learning, broadly defined as practice-based improvement, characterizes many early visual tasks, such as discrimination of complex gratings (Fiorentini and Berardi 1980), hyperacuity (McKee and Westheimer 1978; Poggio et al. 1992), stereopsis (Ramachandran and Braddick 1973; Fendick and Westheimer 1983), discrimination of texture (Karni and Sagi 1991), direction of motion (Ball and Sekuler 1987), or line orientation (Vogels and Orban 1985). A critical issue common to all these studies is the time course over which learning occurs, ranging from minutes (fast learning) to weeks (slow learning). It appears that the fast learning is binocular, indicative of involvement of neural circuitry in which full binocular integration has occurred, while the neural substrate of slow learning may contain a mixture of binocular and monocular neurons as found in the first stage of cortical processing (primary visual cortex, or V1). These attempts to address where in the pathway learning may occur are particularly interesting in the context of recent important evidence for plasticity in the adult early visual system (Frégnac et al. 1988; Gilbert and Wiesel 1992), which provides a significant paradigm shift from the previously prevailing dogma that within early sensory pro-

cessing the properties of the neural mechanisms are subject to experience only during development but are fixed in adulthood (Hubel and Wiesel 1962). This paradigm shift has motivated studies of learning and plasticity at the neuronal level, and the combination of simultaneous neuronal and behavioral investigations (Newsome et al. 1989a; Zohary et al. 1994). For example, Zohary et al. (1994) have documented cortical changes in macaque monkeys trained to discriminate opposite directions of motion in dynamic stochastic random-dot displays (a detailed description of the stimuli is in Newsome and Paré 1988) in which a varying proportion of dots carrying the directional signal was embedded in masking motion noise (Fig. 1a). These type of stimuli illustrate global motion because the extraction of direction cannot be obtained by local computations only, spatial integration of the motion signal over the whole image being required.

Zohary et al. (1994) recorded from neurons in the middle temporal area (MT) while the monkeys were performing a direction discrimination task: within just a few hundred trials the animals demonstrated improvement in their ability to pick out the correct direction of the signal dots, and this concurred with an improvement in the direction specificity of the MT neurons. They suggested that this improvement was an example of perceptual learning. Since the learning transferred from a trained site in the receptive field of an MT neuron to an untrained site, the authors concluded that MT neurons were responsible for learning to discriminate the direction of the motion signal embedded in masking noise, and that learning was not mediated by neurons with a smaller receptive field like those in the early stages of visual pathway, such as in the primary visual cortex (V1).

But what are the mechanisms by which MT neurons mediate learning direction discrimination in these spatially complex stimuli? The goal of the learning process is to improve global direction judgments based on the representation of motion encoded by MT neurons, and thus it must deal with the responses of these units to noise in the motion stimuli. In this paper we suggest two ways of dealing with noise: the first one uses internal template representations corresponding to the directions to be

discriminated, and during learning these templates are altered to *accommodate* the noise. The second one uses a weighted encoding of the representation, and the weights are gradually altered in order eventually to *ignore* units in the representation that are tuned to noise. In Sect. 3 we describe the architecture of the two learning models and in Sect. 4 we present simulated experiments showing that, like the human observers trained on the same task, both models learn in only a few hundred trials to discriminate reliably between opposite directions of motion in stimuli with weak motion signal embedded in masking motion noise. These are the same type of stimuli used by Zohary et al. (1994) and first described in Newsome and Paré (1988). Sect. 2 summarizes the psychophysical learning experiments carried out on human observers. Abbreviated forms of the models and psychophysical results have been presented in (Sundareswaran and Vaina 1995; Vaina et al. 1995).

2 Learning direction discrimination in global motion

In an effort to assess the hypothesis that the time course for learning direction discrimination in global motion is consistent with the fast learning method, we measured performance of naive observers¹ using a training method similar to that proposed by Fiorentini and Berardi (1980) for learning to discriminate complex gratings. The subjects were asked to make judgements about global motion perceived in random-dot kinematograms in which 25% of the dots were displaced in a single direction while the remaining dots were randomly repositioned from frame to frame. The test stimuli consisted of two successively presented frames of different stationary random-dot patterns. Each frame was composed of 100 white dots plotted within an imaginary circular aperture 10 degrees in diameter. In the second frame, 25 of the dots were repeated with a displacement of 6 arc min (either horizontal or vertical displacement, depending on the test). A single frame was 45 ms in duration, and there was no interstimulus interval. The total duration of a trial was 90 ms – too short for subjects to initiate eye movements. Observers were asked to maintain fixation on a small rectangular mark placed at midline 2 degrees to the left or right of the outer margin of the stimulus. Observers' performance improved in the absence of reinforcement (feedback), was fast (200–400 trials), stabilized quickly and was retained for days (Fig. 1b) or even months. These effects of practice were specific to the stimulus location in the visual field (Fig. 1c,d) and to the trained direction of motion (Fig. 1c,e). We showed that the effect of practice transferred from the trained eye to the untrained eye, and that the learning is retinal-specific. Interestingly, not all observers improved their performance. Specifically,

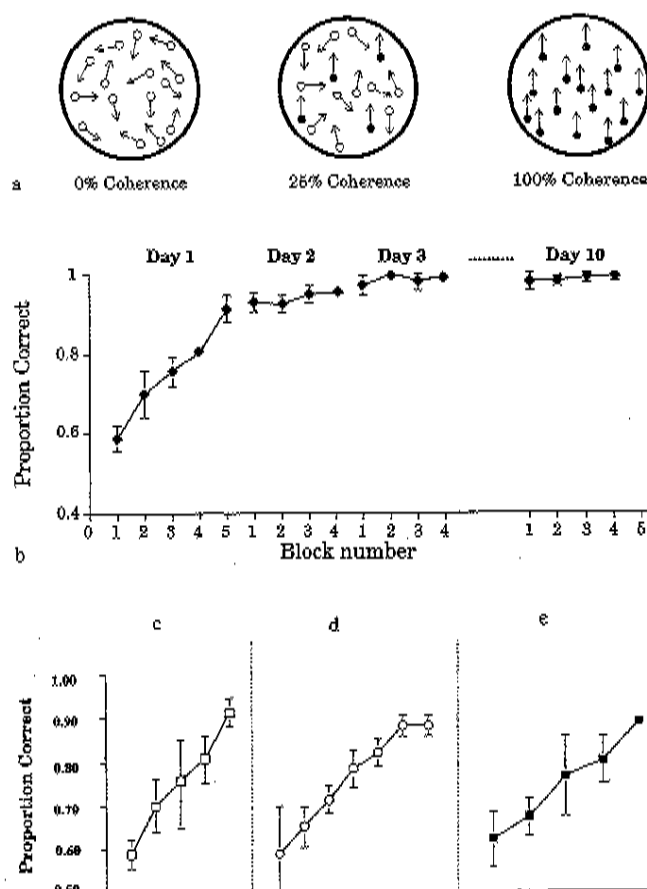


Fig. 1a–e. Psychophysics. a Schematic of the stimulus consisting of random dots each moving either in a coherent motion direction or randomly. The right diagram depicts all dots moving in the coherent motion direction (upward), the left diagram depicts pure random motion, and the middle diagram depicts the stimulus used in the learning experiments, with 25% of the dots moving in the coherent direction and the rest (75%) moving in random directions. b Performance improvement in seven subjects averaged and shown as data points for blocks of 40 trials each of left-right motion with the stimulus at left visual field eccentricity of 2 degrees. c The performance on day 1 of b reproduced. d Performance of the same subjects on day 2 when the stimulus is presented at left visual field eccentricity of 5 degrees, showing no transfer. e Performance on day 2 when the motion direction was changed to up-down

observers starting at chance level were not able to stabilize and improve their performance sufficiently to qualify as learning. Further details and discussion of the psychophysical learning study can be found in Vaina (in preparation).

2.1 Neurobiological basis for the input representation: neurons that discriminate direction in global motion

Global motion is spatially complex such that its direction must be computed from several component motion vectors. Physiological (e.g. Maunsell and van Essen 1983; Newsome et al. 1989b) and behavioral (Newsome and Paré 1988) studies in the macaque found that the MT is the first area in the motion pathway where a large majority of neurons are directionally selective, and thus

¹ Informed consent according to the Boston University institutional review board for research with human subjects was obtained from all subjects.

is a good candidate for encoding global motion. Such a spatially complex task, in which the global direction of the stimulus is extracted from integrating motion information throughout the display, has been shown (Newsome and Paré 1988) to characterize the perceptual abilities of MT neurons. More recently (Britten et al. 1992; Salzman and Newsome 1994) studies involving the simultaneous recording and manipulating of MT neurons while trained monkey were performing a near-threshold discrimination task, revealed that these neurons carry direction signals of sufficient precision to account for psychophysical performance. This result, taken together with the lesion studies, supports the idea that directional signals in the MT contribute directly to the perception of motion, and thus motivates our choice of MT units as input-representation for networks that learn to discriminate opposite direction of motion in the same stimuli used with the human observers, and simulate the psychophysical learning results (Fig. 1).

2.1.1 Relationship between MT neurons and motion input from V1. The analysis of the mean firing rate of the MT neurons and of psychophysical performance in humans and macaque monkeys has revealed an almost linear correlation with changes in the strength of the signal: perfect direction discrimination at 100% coherence, and a linear decrease as the proportion of coherence is decreased (Fig. 7a). The linear correlation-responses of MT neurons to the stochastic random-dot displays described in this section (Fig. 1) can be obtained by linear pooling of local motion filter inputs (Downing and Movshon 1989).

Thus, on the basis of the assumption that learning global motion is mediated by MT neurons, a natural representation of the input to the learning mechanism is a summation of response of the V1 neurons. Units performing summation of local (V1) unit responses, denoted σ units, have response of the form

$$x_i = \sum_{j=1}^n e^{-(1/2\sigma_i^2)(\theta_j - \theta_0)^2} \quad (1)$$

where σ_i is the standard deviation of the tuning curve, and information from n local units is taken into account. The angle θ_j is the preferred motion direction for the unit, and θ_0 the actual motion direction of the i th point in the visual field. Thus, the MT-like unit with a certain preferred direction sums over local units tuned to the same preferred direction.

2.1.2 How may MT neurons compute global motion? Cells in the MT are broadly tuned to the direction of visual motion. Dot patterns moving in different directions generate discharge patterns that when mapped to motion directions result in tuning patterns similar to those described by Georgeopolos et al. (1986) in the motor cortex and by Lehky and Sejnowski (1990) in the visual system as *population coding*. Population coding theories assume that distributed patterns of activity in neuronal populations underlie perceptual behavior, and that correspondingly the neurons will show broadly graded responses. It

is possible that responses of MT neurons can be represented as an instance of population coding and, as such, a few directionally tuned neurons will be used to describe all the directions in the global motion stimulus.

3 Two computational schemes for learning global motion direction

In this paper we focus on modeling the learning rather than on computing the representation of motion prior to learning (this will be addressed in a future study). Here we assume that motion information is available to the learning models in the form of a representation consisting of responses of a collection of directionally tuned MT-like units. A directionally tuned unit is characterized by a response function which has a high value for a certain preferred direction, and decays with angular separation of directions from this preferred direction. Several current computational models have used velocity-tuned filters with similar response properties for computing image velocity (Adelson and Bergen 1985; Heger 1987; Fleet and Jepson 1989).

3.1 Learning to accommodate

We suggest that there is an internal template for global motion in any given direction. In the representation proposed above, the template for a global direction corresponds to the collection of responses from direction-sensitive units to (noise-free) motion in that global direction. Since such templates correspond to a correlation of 100% (i.e., all the dots are moving in the same direction), it is a trivial task to discriminate between leftward and rightward motion by finding the better match among the internal templates for the motion measurement. If we use the same templates for lower values of correlation, direction judgements can be expected to be worse; the performance will deteriorate with lowering correlation, and will eventually reach chance levels (i.e., the judgements are completely random).

The learning to accommodate model gradually alters initial internal templates of global motion. Each initial template is a vector $[x_1, x_2, \dots, x_n]$ corresponding to responses x_i of directionally tuned, MT-like units to noise-free motion in a specific direction. In the presence of noise, however, the responses x_i will depend on the signal-to-noise ratio (SNR), and for a given SNR the response vectors to several inputs form a cluster. The learning process attempts to estimate the center of the clusters corresponding to the left and right noisy global motion, starting from the initial templates as the hypothesized centers. In other words, starting with 'clean' templates (corresponding to zero noise) of the global motion, the method learns to accommodate the noise in the input by gradually altering the templates.

The model is a combination of HyperBF-like functions (Poggio and Girosi 1990) and clustering. We use gaussians with mean at the cluster centers, and 'move' the cluster centers by a learning algorithm. The model is presented schematically in Fig. 3a. The input units x_i

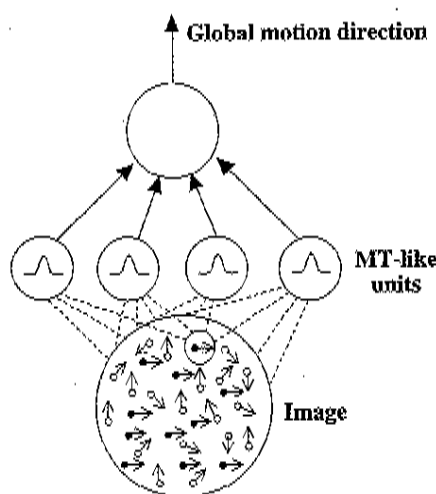


Fig. 2. Generic schematic of the modeling. Units with characteristics similar to those of MT neurons are used in the simulations; these units integrate motion information over a large area by summing over local responses of tuned units. The responses of the MT-like units are used to determine global motion direction and to learn the task.

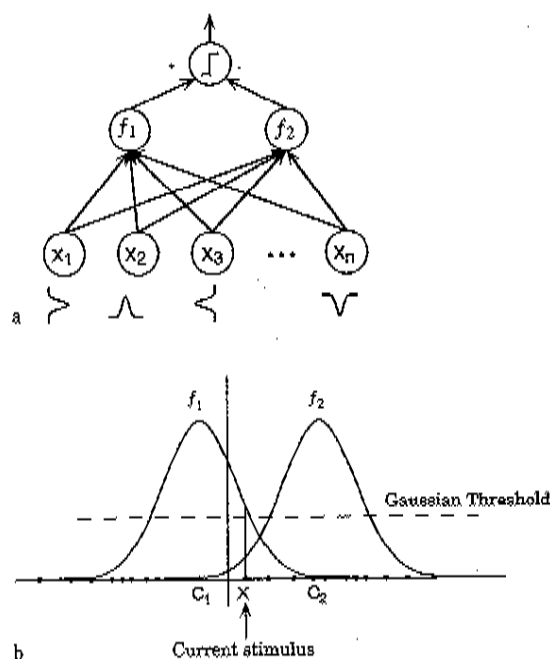


Fig. 3. a Architecture of the learning to accommodate model: inputs x_i represent responses of MT-like units; cluster gaussian functions f_i measure the closeness of the response vector to the corresponding cluster center. Decision rule for the winning cluster is shown schematically as a non-linear thresholding operation on the cluster gaussian responses. b One-dimensional schematic of the classification: winning cluster for a given input X is the one with a higher response above a threshold; if the higher response is below the threshold, the winner is chosen randomly. The cluster centers are marked C_1 and C_2 .

correspond to the responses of directionally tuned neurons. Hidden units f_i correspond to the two cluster gaussians, and the output is a linear combination of the gaussian outputs. A cluster gaussian response is

a measure of how close the current input is to the corresponding cluster center:

$$e^{-(\vec{x} - \vec{c}_i)^T \sigma_i^{-2} (\vec{x} - \vec{c}_i)} \quad (2)$$

where \vec{x} is the current input, \vec{c}_i is the current estimate of the center of cluster i , and $\sigma_i^2 I$ is the covariance matrix of cluster gaussian. The cluster gaussian with the largest response 'wins'. This is shown pictorially for an one-dimensional example in Fig. 3b.

Learning occurs by improvement in the estimate of the cluster centers. The following learning rule (e.g., Moody and Darken 1989; Hertz et al. 1991) modifies the current estimate of the center of the 'winning' cluster:

$$\vec{c}_w^{(t+1)} = \vec{c}_w^{(t)} + \eta \times (\vec{x}^{(t)} - \vec{c}_w^{(t)}) \quad (3)$$

This rule moves the center towards the new data vector \vec{x} that has been judged to belong to the w th cluster. Using this rule, a reliable estimate of the cluster center is obtained after a sufficient number of input presentations. The parameter η controls the learning rate.

In summary, *learning to accommodate* consists of the following steps:

1. Initialize cluster centers to templates corresponding to 100% correlated global motion in the directions to be discriminated.
2. For each trial, repeat the following steps:
 - determine the largest $G_i = \exp(\vec{x} - \vec{c}_i)$; let this be G_w ; if G_w is above a chosen threshold g , (henceforth, gaussian threshold), declare w to be the winning cluster; otherwise choose a cluster w randomly;
 - move the center \vec{c}_w using the rule in (3), and update the standard deviation σ_w ;
 - interpolate using the cluster gaussians as HyperBFs, if desired.

In the above, \vec{c}_i is the i th cluster, and \vec{x} is the current input vector. The use of the gaussian threshold is explained below.

There are several reasons why we use cluster gaussians to determine cluster membership, instead of simply using a distance measure. First, the standard deviation of the gaussian (σ_w) can be *learned*; this will be useful for situations where the noise depends on the signal, i.e., where certain signals are more noise-prone, and have correspondingly 'broader' clusters. Second, the gaussians may be used as basis functions to perform learning as an interpolation scheme, exactly as in HyperBF learning (Poggio and Girosi 1990). Third, the gaussians are useful in modeling the gray area in classification (where one is unable to judge class membership, and makes random judgments) by thresholding the response of the winning cluster gaussian.

3.2 Learning to ignore

In our (MT-like) representation, the global motion direction is encoded as the responses of several direction-sensitive units. This information can be 'decoded' by a weighted combination of the responses. Figure 4 provides a schematic of the architecture of a network that performs the decoding. Let the response of the units

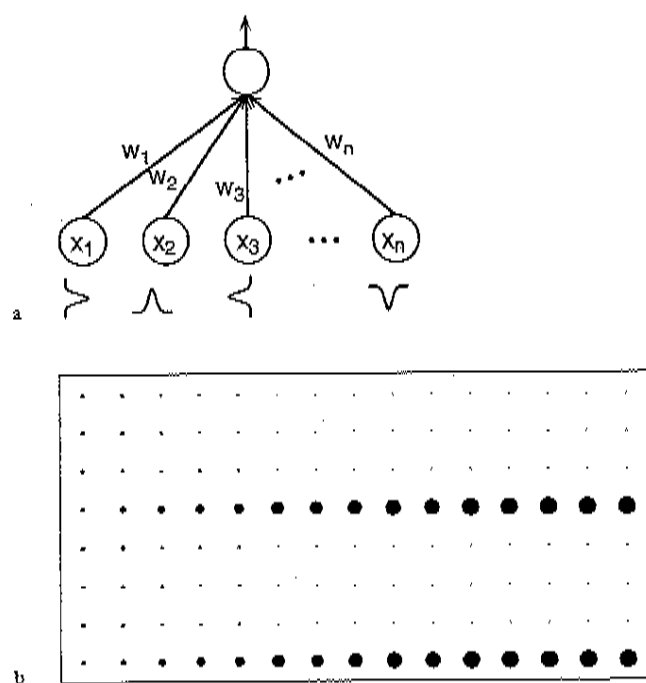


Fig. 4. a Architecture of the *learning to ignore* model: global motion direction is computed as a weighted average of the tuning directions. b Schematic showing changing weights in the model in a typical simulation. Weights of eight directions are depicted as circles (radius proportional to weight value) in columns: the leftmost column corresponds to the weights before training, and each successive column depicts weights after successive blocks of training

be x_1, x_2, \dots, x_n . Let the corresponding weights be w_1, w_2, \dots, w_n . The global velocity is decided as a weighted combination of the individual unit preferred velocities; each weight is a product of adjustable weights w_i and the unit responses x_i . Weighting by the adjustable w_i provides for learning, and weighting by the unit responses x_i assures correct decoding of any instance of the representation. If the neuron preferred velocity is $v_i = [\cos(\theta_i), \sin(\theta_i)]^2$, the output velocity is calculated as

$$v_o = \sum_i w_i x_i v_i \quad (4)$$

or, equivalently, the global motion direction is calculated as

$$\theta_o = \tan^{-1} \left(\frac{\sum w_i x_i \sin \theta_i}{\sum w_i x_i \cos \theta_i} \right)$$

Judgment of the global motion direction is based on θ_o . For instance, the global direction is judged to be 'right' if $+\theta_{thr} > \theta_o > -\theta_{thr}$, and to 'left' if $\pi + \theta_{thr} > \theta_o > \pi - \theta_{thr}$, where θ_{thr} determines an annulus of tolerance within which the computed direction is expected to fall. The use of θ_{thr} in this model is analogous to thresholding cluster gaussian response in the *learning to accommodate* model.

Clearly, the output direction will be strongly biased towards the preferred direction of a unit with a high response and a high connection strength (weight). If the task is to discriminate between leftward and rightward motions, units responding to motion in other directions contribute 'noise'³ to the final computation. The goal of the *learning to ignore* model is to suppress the contribution of these noise neurons by reducing the corresponding weights.

Since prior to learning the model has no bias for one direction over another, all the weights are equal. During learning, the weights w_i are altered by a learning rule. We suggest two learning rules based on two different notions of what information is used to alter the weights: first, an *exposure-based rule* in which units that are consistently active increase their weights, and second, a *self-supervised rule*, in which the model uses its own prediction of global motion direction to change the weights. In both cases, to prevent uncontrolled growth, the weights are renormalized to maintain a unit weight vector (or to maintain the sum of the weights to be unity; in the simulations, there was no qualitative difference in performance between the two choices of renormalization).

3.2.1 Exposure-based rule. In this scheme, the weight corresponding to a unit is incremented by an amount proportional to the current weight. Only units whose response values are above a certain threshold are allowed to increase their weights; this learning rule favors units that are often active. This is consistent with the Hebb rule because a consistently active input must be contributing to the output, and its connection to the output should be strengthened. The learning rule is as follows:

$$w_i \leftarrow w_i + \eta w_i, \quad \text{if } x_i > r_i \quad (5)$$

where r_i is a threshold, and η is a small fraction that controls the learning rate. The same learning rule has been used in learning hyperacuity by Weiss et al. (1993).

It is interesting to note that the model need not be performing a task in order to learn, since the learning rule merely depends on a unit being active, no matter for what reason. Neurons with similar properties have been reported by Zohary et al. (1994), who found neurons in the macaque MT that improved in sensitivity even though the monkey was merely fixating while being exposed to global motion stimuli like the ones used in our study.

3.2.2 Self-supervised rule. In the self-supervised learning mechanism, an internal feedback signal is derived from the computed global motion direction. The weight corresponding to a unit is increased by an amount proportional to the product of the current weight and a measure of the agreement between the computed global motion direction and the preferred direction of the unit:

$$w_i \leftarrow w_i + \eta w_i e^{-(\theta_i - \theta_o)^2 / 2\sigma_i^2}, \quad \text{if } x_i > r_i \quad (6)$$

² We have assumed unit speed, for simplicity.

³ Clearly, the choice of the task influences the notion of what is considered noise.

The model uses its own estimate of the global motion direction to generate an internal reinforcement that controls the learning. For the purpose of learning, the model requires only measures of the discrepancy between the model output and the tuning parameters of input units. Such measures may be available even without a precise knowledge of the output.

This rule also has a Hebbian flavor because units whose weights are increased the most are those that contribute the most to the output. Equivalently, a unit whose response is not strong enough to sway the output closer to its preferred direction is forced to contribute less to the output by a reduction of the corresponding weight (the reduction is not explicit: it occurs due to the renormalization).

In summary, *learning to ignore* consists of the following steps:

1. Initialize all the weights to equal values.
2. For each input presentation:
 - determine the global motion direction as the weighted combination of unit preferred directions;
 - if computed global motion direction is within θ_{thr} of a candidate global motion direction (left or right, for example), declare that candidate direction to be the global motion direction, else choose randomly;
 - alter the weights using either the rule in (5) or the one in (6), and
 - renormalize the weights.

4 Experiments

The purpose of our simulations was primarily to replicate psychophysical performance. In the simulations, the input representation was generated in the following manner. A collection of eight MT-like units, with preferred directions uniformly spread over 2π (four cardinal directions and four oblique directions) was used. For simplicity, speed-selectivity was ignored (unit speed was assumed). The units integrated information from the whole visual field (with no variation due to eccentricity), using (1). For each *trial*, the coherent motion direction was randomly decided. Motion of 40 random dots was simulated, with 10 dots (25% correlation) moving coherently (left or right) and the remaining 30 moving randomly. Trials were grouped into blocks of 50. For the experiments reported here, we simulated the responses of neurons based on their directional tuning. By comparing the decision of a model with the known correct response, each trial outcome was labeled as correct or wrong. The percentage of correct responses was used as a measure of performance (as in the psychophysics).

4.1 Learning to accommodate: simulations

Simulations were performed to study the learning by moving the centers (3). If the representation vector for a trial was close enough to one of the cluster centers, the

vector was assigned to that cluster. 'Closeness' was measured by a measure of the distance to the cluster center; the measure used was the cluster gaussian function value. If a representation vector could not be assigned to either cluster (i.e., both cluster gaussian values for this vector were below a threshold), then the vector was randomly assigned to one of the clusters (this corresponds to the *forced choice* paradigm in psychophysics). We believe that the use of thresholding results in a more realistic model of human decision-making than the case where a sharp decision is made in favor of the nearest cluster.

Typical learning curves, averaged over the performance in 10 simulation runs, are shown in Fig. 5. For the simulation, an η value [see (3)] of 0.0075 was used. Use of this value resulted in performance improvement comparable to that in psychophysics. The cluster gaussians had a fixed σ values of 0.5 (nominal variations in this value did not have any impact on the qualitative performance).

4.2 Learning to ignore: simulations

Global motion direction was computed from the representation vector by the weighted averaging in (4). A simple decision rule would choose the coherent motion directions closest to the computed global motion direction. Our decision rule chooses the coherent motion direction that is closer and is within a certain threshold angular distance (10 degree range), as described in Sect. 3.2. If the global motion direction was within a 10 degree range around the 0 degree direction (right), the model judged the motion to be rightward. If the global motion direction was within a 10 degree range around the leftward direction, the model

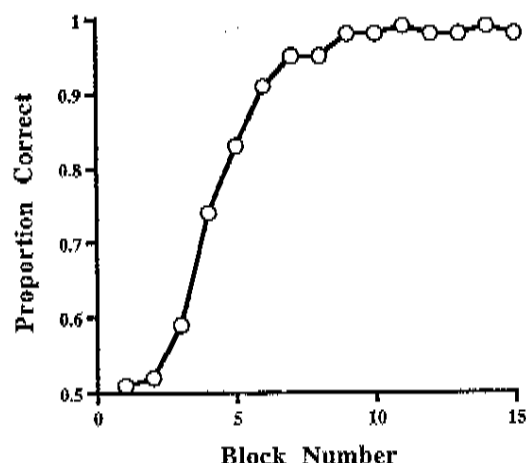


Fig. 5. Performance improvement over time for the *learning to accommodate* model in the simulations. Data have been averaged over 10 simulation runs. σ units tuned to four cardinal and four oblique directions were used. Parameter values: $\eta = 0.0075$, $g_i = 0.8$, $\sigma_h = \pi/8$, $\sigma_s = 0.5$

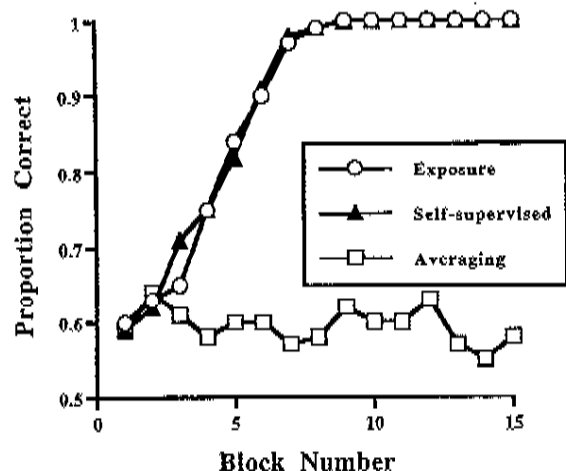


Fig. 6. Performance improvement over time for the *learning to ignore* model in the simulations. Data have been averaged over 10 simulation runs. σ units tuned to four cardinal and four oblique directions were used. Parameter values: $\eta = 0.015$, $\theta_{thr} = 5$ degrees, $\sigma_n = \pi/8$. Results are shown for exposure-based learning (open circles), self-supervised learning (triangles) and a non-learning averaging mechanism (open squares). See text for details

judged the motion to be leftward. In all other cases, the model's judgment was randomized to choose left and right directions with equal probability.

Typical learning curves, averaged over 10 simulations, are shown in Fig. 6. A learning rate (η) of 0.015 and an angular threshold θ_{thr} of 5 degrees were used. The response threshold, r_t , was set to be $0.9 \max(x_i)$. It was found that in all the simulations, the model learned at a rate comparable to human subjects, irrespective of whether the *exposure-based* learning rule or the *self-supervised* learning rule was used.

4.3 Justification of architectures

To verify that the architecture of the models was appropriate for the chosen task, we disabled the learning, and studied the performance of the models for varying proportions of correlated dots (signal). For 100% correlation (no noise), the performance was perfect (100% correct responses), and by decreasing the proportion of correlated dots, the performance level was reduced.

For the *learning to accommodate* architecture (with cluster centers set to positions corresponding to pure left and right motion—100% correlation), we found that the model's performance (Fig. 7b) was very different from that of the monkeys (Fig. 7a). Figure 7b shows the performance of the model for different values of the threshold on the gaussians as the percentage coherence is increased. Clearly, the larger the threshold, the poorer the performance (recall that below threshold, random judgments are made), which explains the shift in the curve for the different threshold values. However, unlike the performance of monkeys or humans on this task, in the model, all the curves stay at chance level for a range of coherence values, and then abruptly climb up to peak performance. This is not surprising, since for low values of coherence

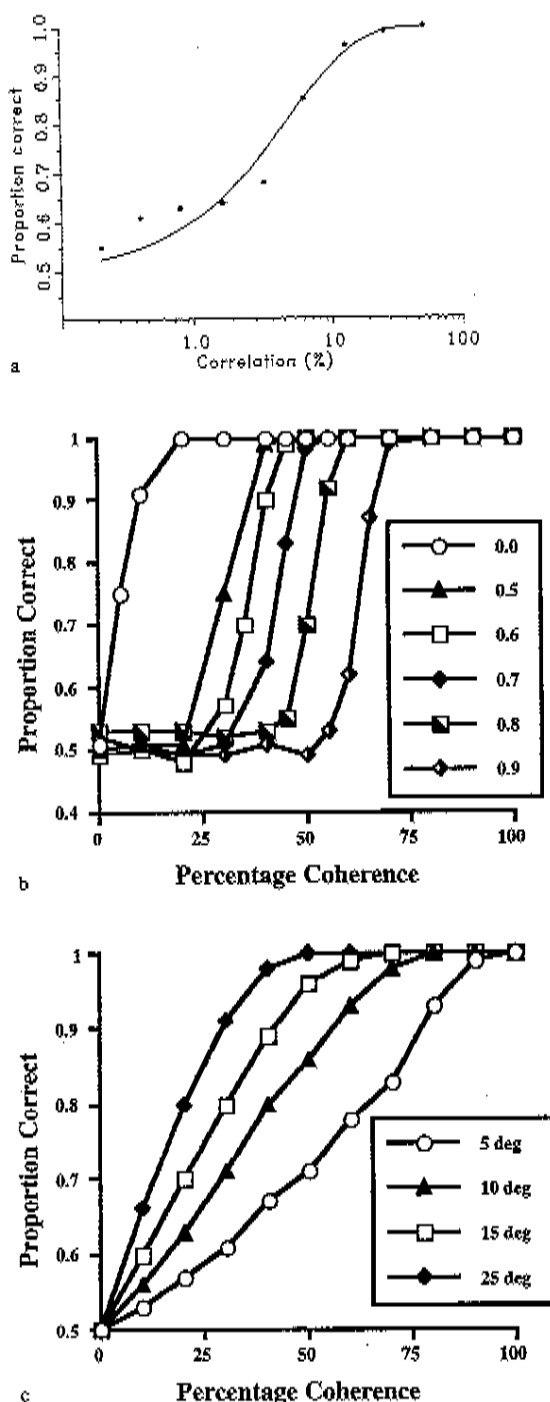


Fig. 7. **a** Psychophysical performance of monkeys for increasing stimulus correlation (reproduced from Newsome et al. 1989b). **b** *Learning to accommodate*: curves for different values of gaussian threshold g_1 , showing performance of the model to increasing signal in the input. The proportion of inputs for which a random decision is made is plotted against the percentage coherence for various values of the gaussian threshold. **c** *Learning to ignore*: curves for different values of the angular threshold θ_{thr} for increasing percentage coherence

the cluster is so 'spread out' that the gaussian threshold cannot be reached (the model has fixed centers corresponding to 100% correlation), and thus the performance is at chance level. However, beyond a certain level of

correlation the cluster is sufficiently close to the fixed center, and the performance jumps to 100% correct. This is what we see in Fig. 7b. The input representation vectors lie largely in the below-threshold region for low values of coherence, and fall within the above-threshold region for high values of correlation, with a rather abrupt transition. This non-smooth behavior suggests that the learning to accommodate model does not have the architecture appropriate to explain human and monkey data.

For the *learning to ignore* architecture, with the learning disabled, the correlation varied and all the weights set equal, the model's behavior conforms to the expected drop in performance with decreasing correlation (Fig. 7c). At 100% correlation, the response of the unit tuned to the coherent motion direction is significantly higher than the response of the other units, resulting in correct decisions. However, when the correlation is reduced, the responses become more evenly distributed, and the direction judgments are more erroneous. The various curves in Fig. 7c show the performance for different values of the angular threshold. As expected, when the angular threshold is large (i.e., a greater range of angles are accepted as being correct), the performance is better.

We also investigated an alternative model of global integration, and we comment on it in Appendix A. In Appendix B, we report miscellaneous experimental observations that shed more light on the behavior of the models.

5 Discussion

5.1 Other learning models that simulate psychophysical performance

While psychophysical and neurophysiological studies of learning visual tasks consistently indicate the involvement of early portions of the visual system whose neuronal properties and functional cortical architecture are modifiable by experience, they do not address explicitly what is being learned, and what may be the plausible computations that underlie learning of a perceptual task. Recently, Poggio and Girosi's (1990) theoretical proposal, that learning from examples may be viewed as synthesizing task-related modules at the level of the cortex, has provided the basis for a series of important studies which explicitly simulate perceptual learning obtained in human observers in several tasks (Poggio et al. 1992; Weiss et al. 1993). Briefly stated, Poggio and Girosi's (1990) hypothesis is that synthesizing a template from examples for a specific task can be viewed as a problem of approximating a multivariate function from sparse data. They showed that the solution to the approximation problem can be expressed in terms of a class of multilayer networks that they named HyperBF functions, which are a type of generalized radial basis functions (GRBFs). HyperBF networks with gaussian basis functions (gaussian HyperBF) are an efficient and neurobiologically plausible model for this computation (Poggio 1990). That

the cortex might actually be constructing such templates through learning is consistent with psychophysical findings of the specificity of learning to the stimulus attributes. Furthermore, Poggio (1990) argued that HyperBF networks with gaussian basis functions can be implemented in terms of biologically plausible mechanisms and circuitry mediating the learning of perceptual tasks. Indeed, brain architecture can easily implement radial basis functions (Moody and Darken 1989; Poggio 1990) as they are radially symmetric and die off exponentially. The gaussian HyperBF is also consistent with the experiment-dependent plasticity of the neuronal and functional architecture, since both the centers of the radial functions and the weights in the norm are updated during learning. The centers represent templates, and updating (moving) the centers is equivalent to modifying the corresponding template (Poggio 1990). Finding the optimal weights corresponds to task-dependent dimensionality reduction, that is, through learning some features become more salient than others. This approach is exemplified by a series of simulated psychophysical experiments in which a HyperBF network was trained in the supervised mode to improve on several hyperacuity tasks. For example, the network, like the human observers, upon briefly being presented with two vertical bars arranged one under the other with a small spatial offset, had to determine whether the upper bar was to the left or to the right of the lower one. Units with gaussian receptive fields were stimulated by input images of the line segments, and by comparing the prediction of the network to known offset direction, the centers of the gaussians and the weights that combine the gaussian responses were adjusted to minimize the error in the output. This network learned to improve visual hyperacuity at a rate comparable to humans (Poggio et al. 1992). A somewhat different approach to modeling learning hyperacuity was taken by Weiss et al. (1993). They used unsupervised learning rules in a HyperBF network with the input set to response characteristics of orientation-selective units and the output to a weighted combination of the basis functions. In this case also the model learned to perform well in the vernier hyperacuity task.

The two models proposed here have close similarities with the approaches of Poggio et al. (1992) and Weiss et al. (1993). The cluster gaussians in the *learning to accommodate* model are similar to HyperBFs (Poggio and Girosi 1990), since in both learning is mediated by moving the centers of radial functions. However, in contrast to HyperBFs, we do not interpolate available data, and thus in our model there are no weights to be modified. Another major difference is that Poggio et al. use supervised learning while we use unsupervised learning. In this respect, our model is closer to that of Moody and Darken (1989), who use cluster gaussians with movable centers. But again, their goal is to obtain interpolation of the input data, and to examine a hybrid learning scheme to move the centers and to alter weights in their network; we focus on classification and on the role of clustering to accommodate noise in the data. The *learning to ignore* model is somewhat similar to the model proposed by

Weiss et al. (1993). Both use a weighted combination of tuned detector responses (motion in the case of the former, and orientation in the latter), and learn by changing the weight values. The exposure-based learning rule is also examined in both models. However, in Weiss et al. (1993) noise is due to the noise neurons, while in our modeling noise is in the input (motion stimulus).

5.2 Can global motion representation be learnt?

Two alternatives have been proposed to compute global motion direction. The first is a *Winner Take All (WTA)* mechanism, and the second is a simple *averaging* mechanism (Salzman and Newsome 1994).

In the WTA method, global motion direction is the preferred direction of the unit with the largest response. This mechanism is necessarily limited, because it will obtain reliable results only if the noise content is relatively low. This does not, however, rule out WTA as a candidate. Salzman and Newsome (1994) present evidence that WTA is more likely than averaging in a global motion direction task. By electrical simulation of directionally tuned MT neurons in a monkey, they verified that motion direction judgment of the monkey was biased towards the most active neuron. Interestingly, our *learning to accommodate* is a clustering approach which is an implementation of WTA (Hertz et al. 1992 discuss clustering as WTA). After learning, each cluster corresponds to a possible global motion direction, and the 'winner' (the one with the largest cluster gaussian response above threshold) 'takes' the current input. We propose that *learning to accommodate* provides a mechanism to learn the WTA method to perform the global motion direction task.

In the averaging method, global motion direction is the average direction of optical flow vectors. To use this alternative, two issues need to be addressed: computation of the optical flow field, and the performance of the averaging mechanism. While there is substantial evidence for directionally tuned neurons in the cortex, there is no direct evidence for the computation of optical flow vectors. Also, methods which have been proposed to compute optical flow from directionally tuned units (e.g., Heeger 1987; Grzywacz and Yuille 1990) lack reliability and robustness. Other methods for computing optical flow (image intensity gradient-based methods and feature correspondence-based methods) exhibit similar shortcomings, and their biological plausibility is questionable. Even assuming that these criticisms are somehow countered, our simulations⁴ show that the averaging mechanism performs worse than the learning methods we presented (Fig. 6). However, *learning to ignore* is an averaging mechanism, and it performs better (than averaging optical flow vectors) because of weighted averaging and alteration of the weights to ignore noise. We suggest that *learning to ignore* is a method that learns an adaptive averaging mechanism.

6 Conclusions

Global motion perception is a critical aspect for the vast majority of our daily activities. When we move through the environment by driving, biking or walking, we use motion within large areas of the visual field to determine our own motion, motion of objects, presence of obstacles in our path, time taken by a moving object to arrive close to us, and the shape of objects. Several computational studies have attempted to characterize mechanisms of global motion perception (e.g., Prazdny 1980; Rieger and Lawton 1985; Heeger and Jepson 1992; Hummel and Sundaeswaran 1993). The methods proposed in these studies can integrate motion information from a large number of points in the visual field to determine the motion of the observer. However, good performance of these methods depends on having little or no noise. In real life, vision-based systems have to deal with various contingencies (fog, rain, smoke, snow, low light, etc.) that considerably degrade visibility. Under such low-visibility conditions, noise-sensitive methods can be expected to fail due to their sensitivity to perturbations in the input.

An interesting question is whether poor performance in motion judgments under low-visibility conditions can be improved by training. In our experiments involving judgment of direction of motion embedded in masking motion noise, we found that performance of subjects improves very rapidly. This performance improvement is sustained, and does not degrade over a long period of time. These psychophysical results suggested the interesting possibility of designing adaptive computational models that can improve on motion perception by training. Such models can be expected to adjust rapidly to cope with changing conditions. Motivated by this potential, and to understand the neural mechanisms behind performance improvement in the psychophysical tasks, we developed two models for fast learning of global motion direction. In this paper, we presented these models and simulation results.

The masking motion noise used in our experiments is only one possible rendering of low-visibility conditions. In further experiments, we plan to explore more realistic presentations of low-visibility conditions. Other experiments currently in progress examine more complex motion scenarios (e.g., expansion and contraction) that are important in three-dimensional motion perception.

In conclusion, our experiments have suggested that there are interesting aspects of plasticity in the motion pathway of the adult human cortex, and in this paper we presented two computational models of this plasticity. The model that learns to perform a weighted averaging seems the more appropriate model of human performance.

Appendices

A Alternative model of MT neurons

Since the stimuli used in the learning experiments are noisy, and because the 'noise' dots carry a mixed motion

⁴The same angular threshold values (see Sect. 4.2) were used in the simulations of averaging as well as *learning to ignore*.

signal and stimulate most MT neurons, it is likely that in reality we have to deal with a mixture of excitation and inhibition to the MT neuron. Inhibition may reduce the effective gain of the excitatory inputs, and this could be represented by encoding MT as a product of V1 responses. The response of a product (π) unit, is defined by

$$x_i = e^{-(1/2\sigma^2 \sum_{k=1}^n (\theta_i - \theta_k)^2)} \quad (A1)$$

For a given preferred direction θ_i , the response is the product of local responses of gaussian units tuned to the same preferred direction. Both σ and π units model adequately the directionally tuned response of MT neurons. The π units do not have such a direct physiological correlate as the σ units; Durbin and Rumelhart (1989) suggested a neurobiological interpretation for π units by hypothesizing a logarithmic transformation at the presynaptic stage followed by a qualitatively exponential response function at the postsynaptic stage. We compare the performance of the two types of units for both *learning to accommodate* and *learning to ignore* models.

Since the product can decay very rapidly as a function of the number of units involved in the product, the spread (σ_k) of the tuning function has to be large compared with those for the σ units. In other words, to make π units responsive, the local (V1) motion units have to respond less selectively to direction. This reduced selectivity to direction renders the population coding 'uniform', especially in the presence of noise. That is, the σ unit responding maximally corresponds to the global motion direction almost always, whereas the π unit responding maximally signals global motion direction to a much lesser degree. Another way to interpret this is that the representation vector is elongated along one of the axes for the σ units but not for the π units and, thus, a winner-takeall mechanism would profit from σ units but not from π units.

We observed the following effects of the uniform representation based on the π units:

- the π units required a high value of the tuning function spread (σ_k); otherwise, many local units have negligible response, pulling down the overall response;
- the clusters of the representation vectors had a greater spread than in the case of σ units, which forced the choice of a larger cluster spread [σ_s , see (2)];
- in learning using the exposure-based rule (5), the smooth representation results in favoring units that are not tuned to the directions being discriminated. This makes it hard to learn using the exposure-based rule, since the 'noise' directions are also favored quite frequently; and
- in general, the learning based on σ unit representation was much more robust to changes in parameter values, and to minor alterations in the learning rules.

B Experimental observations

The results of our simulations suggest that the parameters of the models play a systematic role in the performance. In general, the parameters may be varied over

a reasonably wide range with smooth effects on the performance. The exception is σ_k , the spread of the tuning function. The larger the σ_k , the more uniform the representation (a representation is considered uniform if the variation among its constituent elements is not large). The uniformity of the representation influences learning by the exposure-based rule. This is easily explained. For the exposure-based learning rule to induce the proper weight structure, the favored units (i.e., whose weights are increased) must be the ones pertinent to improve performance in the direction discrimination task (for example, the left and right direction-tuned units, for a left-right direction discrimination task). A uniform representation favors nearby units quite frequently, resulting in poor exposure-based learning. In particular, as discussed in Appendix A, the π units require a large σ_k , which leads to poor and unstable exposure-based learning.

The parameters specific to the learning to accommodate model are the spread (σ_s) of the cluster gaussians, and the gaussian threshold (g_i). The spread of the cluster gaussians, σ_s , and the threshold g_i together account for the initial performance. A sufficiently low σ_s and a sufficiently large g_i are needed to mimic the poor performance exhibited by the human subjects in the first block of the learning experiments.

The parameters specific to the *learning to ignore* model are angular threshold θ_{thr} and threshold r_i . The angular threshold θ_{thr} plays the role of σ_s and g_i discussed above, that is, to tune the initial performance. The role of r_i is, however, more complex. The thresholding operation may be omitted for the σ units but not for the π units due to the uniformity of the π representation. With thresholding, the learning is more robust (less susceptible to changes in other parameters and signal content in the input), since by thresholding only the highly active units (responding to signal) have the possibility of increasing their weights.

For the *learning to ignore* model, where σ units are used, the exact form of the learning rules is not very critical. Learning occurs with any of the following exposure-based rules, which differ from (5) by replacing w_i with either x_i or $w_i x_i$:

$$w_i \leftarrow w_i + \eta x_i, \quad \text{if } x_i > r_i,$$

or

$$w_i \leftarrow w_i + \eta w_i x_i, \quad \text{if } x_i > r_i$$

In the latter case, η must have a higher value to compensate for the reduction due to the product $w_i x_i$. The same changes to the self-supervised learning rule (6) do not alter learning performance. From the empirical evidence, we conclude that *learning to ignore* is a robust form of learning.

Acknowledgements. This research was conducted at the Intelligent Systems Laboratory of Boston University, College of Engineering. L.M.V. and V.S. were supported in part by grants from the Office of Naval Research (# N00014-93-1-0381) and the National Institute of Health (EY 2R01-07861) to L.M.V. V.S. was in part supported by the Boston University College of Engineering Dean's postdoctoral fellowship. We gratefully acknowledge additional financial support for V.S. from the Dean's special research fund.

References

- Adelson E, Bergen JR (1986) The extraction of spatio-temporal energy in human and machine vision. *IEEE Workshop on Motion*, pp 151-155
- Ball K, Sekuler R (1987) Direction-specific improvement in motion discrimination. *Vision Res* 27:953-965
- Britten KH, Shadlen NM, Newsome WT (1992) The analysis of visual motion: a comparison of neuronal and psychophysical performance. *J Neurosci* 12:4745-4765
- Downing CJ, Movshon JA (1989) Spatial and temporal summation in the detection of motion in stochastic random dot displays. *Invest Ophthalmol Vis Sci [Suppl]* 30:72
- Durbin R, Rumelhart DE (1989) Product units: a computationally powerful and biologically plausible extension to backpropagation networks. *Neural Comput* 1:133-142
- Fendick M, Westheimer G (1983) Effects of practice and the separation of test targets on foveal and perifoveal hyperacuity. *Vision Res* 23:145-150
- Fiorentini A, Berardi N (1980) Perceptual learning specific for orientation and spatial frequency. *Nature* 287:43-44
- Fleet DJ, Jepson AD (1990) Computation of component image velocity from local phase information. *Int J Comput Vis* 5:77-104
- Frégnac Y, Shulz D, Thorpe S, Bienenstock E (1988) A cellular analogue of visual cortical plasticity. *Nature* 333:367-370
- Georgopoulos AP, Schwartz AB, Kettner RE (1986) Neuronal population coding of movement direction. *Science* 233:1416-1419
- Gilbert CD, Wiesel TN (1992) Receptive field dynamics in adult primary visual cortex. *Nature* 356:150-152
- Grzywacz NM, Yuille AL (1990) A model for the estimate of local image velocity by cells in the visual cortex. *Proc R Soc Lond A* 239:129-161
- Heeger DJ (1987) Optical flow using spatiotemporal filters. *Int J Comput Vis* 1:279-302
- Heeger DJ, Jepson AD (1992) Subspace methods for recovering rigid motion. I. Algorithm and implementation. *Int J Comput Vis* 7:95-117
- Hertz J, Krogh A, Palmer RG (1991) Introduction to the theory of neural computation. Addison-Wesley, Reading, Mass
- Hubel D, Wiesel T (1962) Receptive fields, binocular interaction, and functional architecture in the cat's visual cortex. *J Physiol (Lond)* 160:106-154
- Hummel R, Sundaeswaran V (1993) Motion parameter estimation from global flow field data. *IEEE Trans Pattern Analysis Machine Intell* 15:459-476
- Karni A, Sagi D (1991) Where practice makes perfect in texture discrimination: evidence for primary visual cortex plasticity. *Proc Natl Acad Sci USA* 88:4966-4970
- Lehky SR, Sejnowski T (1990) Neural model of stereoacuity and depth interpolation based on a distributed representation of stereo disparity. *J Neurosci* 10:2281-2299
- Maunsell JHR, Essen DC van (1983) Functional properties of neurons in the middle temporal visual area of the macaque monkey. I. Selectivity for stimulus direction, speed, and orientation. *J Neurophysiol* 49:1127-1147
- McKee SP, Westheimer G (1978) Improvement in vernier acuity with practice. *Percept Psychophys* 24:258-262
- Moody J, Darken CJ (1989) Fast learning in networks of locally-tuned processing units. *Neural Comput* 1:281-294
- Newsome WT, Paré EB (1988) A selective impairment of motion perception following lesions of the middle temporal visual area. *J Neurosci* 8:2201-2211
- Newsome WT, Britten KH, Movshon JA (1989a) Neuronal correlates of a perceptual decision. *Nature* 341:52-54
- Newsome WT, Britten KH, Movshon JA, Shadlen NM (1989b) Single neurons and the perception of visual motion. In: Lam DMK, Gilbert CD (eds) *Neural mechanisms of visual perception: proceedings of the retinal research foundation*. Portfolio Publishing, Huntington, NY, pp 171-198
- Poggio T (1990) A theory of how the brain might work. *Cold Spring Harbor Symp Quant Biol* 55:899-910
- Poggio T, Girosi F (1990) Networks for approximation and learning. *Proc IEEE* 78:1481-1497
- Poggio T, Fahle M, Edelman S (1992) Fast perceptual learning in visual hyperacuity. *Science* 256:1018-1021
- Prazdny K (1980) Egomotion and relative depth from optical flow. *Biol Cybern* 36:87-102
- Ramachandran VS, Braddick O (1973) Orientation-specific learning in stereopsis. *Perception* 2:371-376
- Rieger JH, Lawton DT (1985) Processing differential image motion. *J Opt Soc Am A* 2:354
- Salzman CD, Newsome WT (1994) Neural mechanisms for forming a perceptual decision. *Science* 264:231-237
- Sundaeswaran V, Vaina LM (1995) Learning direction in global motion: two classes of psychophysically-motivated models. In: Tesauro G, Touretzky D, Leen T (eds) *Advances in neural information processing systems 7*. 7:917-924, The MIT Press, Cambridge, Mass.
- Vaina LM, Sundaeswaran V, Harris J (1995) Learning to ignore: psychophysics and computational modeling of fast learning of direction in noisy motion stimuli. *Cogn Brain Res* 2:155-163
- Vogels R, Orban GA (1985) The effect of practice on the oblique effect in line orientation judgements. *Vis Res* 25:1679-1687
- Weiss Y, Edelman S, Fahle M (1993) Models of perceptual learning in vernier hyperacuity. *Neural Comput* 5:695-718
- Zohary E, Celebrini S, Britten KH, Newsome WT (1994) Neuronal plasticity that underlies improvement in perceptual performance. *Science* 263:1289-1292