

## Attention as Reward-Driven Optimization of Sensory Processing

**Matthew Chalk**

*matthewjchalk@gmail.com*

*Group for Neural Theory, LNC, DEC, Ecole Supérieure, Paris 75005, France*

**Iain Murray**

*i.murray@ed.ac.uk*

**Peggy Seriès**

*pseries@inf.ed.ac.uk*

*Institute for Adaptive and Neural Computation, School of Informatics, University of Edinburgh, Edinburgh EH8 9AB, U.K.*

**Attention causes diverse changes to visual neuron responses, including alterations in receptive field structure, and firing rates. A common theoretical approach to investigate why sensory neurons behave as they do is based on the efficient coding hypothesis: that sensory processing is optimized toward the statistics of the received input. We extend this approach to account for the influence of task demands, hypothesizing that the brain learns a probabilistic model of both the sensory input and reward received for performing different actions. Attention-dependent changes to neural responses reflect optimization of this internal model to deal with changes in the sensory environment (stimulus statistics) and behavioral demands (reward statistics). We use this framework to construct a simple model of visual processing that is able to replicate a number of attention-dependent changes to the responses of neurons in the midlevel visual cortices. The model is consistent with and provides a normative explanation for recent divisive normalization models of attention (Reynolds & Heeger, 2009).**

### 1 Introduction ---

Attention plays an important role in sensory perception, improving one's perceptual performance at detecting attended stimuli, at the expense of a reduction in performance for other stimuli (Pestilli & Carrasco, 2005). A large body of work has been devoted to identifying the neurophysiological changes underlying attention-dependent changes in perception

---

A supplemental appendix is available online at [http://www.mitpressjournals.org/doi/suppl/10.1162/NECO\\_a.00494](http://www.mitpressjournals.org/doi/suppl/10.1162/NECO_a.00494).

(Reynolds & Chelazzi, 2004). A central finding has been that in the striate and extrastriate visual cortex, the firing rate of neurons tuned toward attended spatial locations or features is increased (Reynolds, Pasternak, & Desimone, 2000). Taken alone, this result appears to paint a simple picture: that attention acts to optimize sensory processing toward attended stimuli by increasing the sensitivity of sensory neurons that are tuned toward these stimuli. However, on closer inspection of the experimental data, it becomes clear that this picture is overly simple. In addition to increasing neural firing rates, visual attention can also suppress responses (Reynolds, Chelazzi, & Desimone, 1999), alter receptive field properties (Womelsdorf, Anton-Erxleben, Pieper, & Treue, 2006), and influence center-surround suppression from a stimulus placed outside the classical receptive field (Sundberg, Mitchell, & Reynolds, 2009). Furthermore, the effects of attention are highly sensitive to the experimental setup, with changes in the sensory stimulus and behavioral task giving rise to qualitatively different attention-dependent changes in neural responses.

Recently several divisive normalization models have been proposed that are able to account for many of the experimentally observed effects of attention in the low to midlevel visual cortices (Reynolds & Heeger, 2009; Lee & Maunsell, 2009; Ghose, 2009). While the details of these models vary, the firing rate of a neuron is generally computed by dividing its feedforward excitatory input by the summed activity of a pool of neurons with similar, but differing, stimulus selectivities. These models explain why attention can facilitate or suppress the response of a given neuron, depending on how it alters the neuron's excitatory input, versus suppression from other neurons. They also provide a potential explanation as to why small changes in the behavioral task can produce qualitatively different types of attentional modulation. For example, if the task requires the animal to attend to a small region of visual space, then the principal effect of attention will be to alter the feedforward input to a neuron that is tuned to this location, giving rise to simple multiplicative changes in its firing rate. Alternatively, if the task requires the animal to direct its attention toward a broader region of space, then attention will alter the activity of neurons tuned to nearby spatial locations also, increasing both the inhibitory and excitatory input to a neuron that is tuned to the center of the attended spatial region. As a result, the neuron will undergo a more complex form of attentional modulation that cannot be explained by a simple multiplicative change in its firing rate.

A limitation of divisive normalization models of attention is that the modulatory effect of attention on the feedforward input to each of the neurons in the network has to be specified explicitly by the modeler rather than being predicted directly from the behavioral task and presented visual stimuli. To avoid this limitation, we need a theory that can explain why, rather than just how, attention alters sensory neural responses as it does.

Several researchers have proposed that attention-dependent changes to sensory neural responses can be understood within a normative Bayesian

framework as a consequence of performing optimal inferences about the state of the world (Dayan & Zemel, 1999; Rao, 2005; Chikkerur, Serre, Tan, & Poggio, 2010; Yu & Dayan, 2005; Yu, Dayan, & Cohen, 2009; Dayan & Solomon, 2010; Whiteley & Sahani, 2012). These models hypothesize that changes to the attentional state of the animal correspond to changes in their prior beliefs about the world, which in turn alter how incoming sensory signals are used to infer which stimuli are present. Recently Chikerrur et al. (2010) showed that given certain assumptions about how probabilistic inference is performed in the brain, increasing one's prior belief that certain (attended) stimuli will be presented produces qualitatively similar changes to neural firing rates as divisive normalization models of attention. However, Chikerrur et al. specified explicitly the attention-dependent changes to the prior without providing a normative explanation for why these changes come about. Indeed a general problem of Bayesian models of attention is that it is often not clear why the animal should alter its prior beliefs, depending on the behavioral task that it is performing. Specifically, in the case where attention is manipulated by changes to the behavioral task (i.e., by manipulating which stimuli are important in determining the action that the animal should perform; Pestilli & Carrasco, 2005; Luck, Chelazzi, Hilliard, & Desimone, 1997), rather than by the presented stimulus statistics (i.e., by manipulating which stimuli are most likely to be presented (Posner, Snyder, & Davidson, 1980; Downing, 1988) there is no clear normative reason that the animal should alter its prior beliefs about which stimuli are most likely to be presented.

Here, we extend previous Bayesian models of attention to account for task-dependent modulation of sensory neural responses. We hypothesize that the nervous system learns an internal model describing how both the sensory input and the reward received for performing different actions are generated by a common set of explanatory causes (Sahani, 2004). Within this framework, the behavioral task will alter visual neuron responses only when there is some mismatch between the organism's internal model of the sensory input and the external environment. We argue that due to the complexity of real-world environments, this is often the case. Faced with such a model mismatch, we propose that attention modulates visual processing in order to improve the organism's predictions of the received reward, at the possible expense of their learning a worse model of the stimulus statistics.

We implement a simple model of visual processing to illustrate how our framework can be used to predict attention-dependent changes to visual neuron responses. For our simulations, we assume a particular type of model mismatch in which the image features that are relevant to the task are smaller than the image features used by the agent to perform the task. In common with previous Bayesian models of attention, we assume that attention alters the internal model in a computationally simple way: varying the prior probability that image features are present while leaving other aspects

of the model unchanged. Given certain assumptions about the form of the internal model and how probability distributions are encoded by the sensory neural population, our model predicts attention-dependent changes to visual neuron responses that are consistent with a number of experimental observations in midlevel regions of the visual cortex, including modulation of contrast response functions, sensory tuning curves, and center-surround interactions. Our model is consistent with and provides a normative explanation for previous divisive normalization models of attention (Reynolds & Heeger, 2009; Lee & Maunsell, 2009; Ghose, 2009).

## 2 Overview of Modeling Approach

---

**2.1 General Framework.** A large body of research is based on the idea that the visual system learns a probabilistic model of natural image statistics, in which a set of hidden causes is assumed to generate received sensory signals (Hyvärinen, 2010). We extend this framework to consider visual processing within the context of a simple task, where a biological agent has to perform actions (motor commands or perceptual judgments) in order to receive a reward. To perform the task, the agent must be able to predict the reward associated with each possible action. We hypothesize that it does this by learning a probabilistic model that describes how both the sensory input and reward received for performing an action are generated by a common set of hidden causes (Sahani, 2004). This internal model is used to infer the hidden causes that generated its received sensory input and, consequently, to predict the reward associated with each action.

In most statistical models of visual processing, the agent's internal model is learned in an unsupervised manner based on the statistics of received sensory signals (Hyvärinen, 2010). We propose that in addition, the internal model is continuously adapted based on received sensory signals and reward in order to optimize performance for the task at hand. As well as influencing behavioral performance, changes in the internal model will also influence perceptual inference, altering the agent's internal representation of sensory stimuli. As a result, the activity of visual neurons will vary dynamically in response to changes in both reward contingencies and presented stimulus statistics. Here, we propose that this task-dependent optimization of the agent's internal model can account for experimentally observed changes in visual neuron responses normally attributed to selective attention.

**2.2 When Do Task Demands Alter Visual Processing?** The responses of visual neurons to a given stimulus can be manipulated by changes in the stimulus statistics (determining which stimuli are expected; often communicated by visual cues) (Posner et al., 1980) or the reward delivered for performing each action (determining which stimuli are deemed relevant to the task) (Pestilli & Carrasco, 2005). In our theoretical framework, the

agent's internal model of the stimulus statistics is coupled to its internal model of reward. Thus, perceptual inference can be altered by changes to both the stimulus and reward statistics. In contrast, previous Bayesian models of visual processing, in which the agent's internal model is learned and adapted based on the stimulus statistics alone, can account only for changes in perception due to changes in the stimulus statistics.

For the agent's internal model of the sensory input statistics to be altered by the reward structure of a task, there must be some mismatch between its internal model and the external environment (i.e., if the internal model is already a perfect description of the world, it cannot be further optimized). We postulate that due to the complexity of real-world environments, this will often be the case. For our simulations, we assume that the image features relevant to the task are more spatially localized than the image features used by the agent to choose which action to perform. Such a model mismatch might occur because the agent tries to learn a simple model of the behavioral task, in which the actions that it should perform depend on a small number of spatially distributed image features. While useful in allowing the agent to quickly learn new tasks, this model structure could result in suboptimal performance in experiments that use very simple or spatially localized stimuli (e.g., orientated gratings or coherent motion).

### 3 Methods

---

We constructed a simple model to illustrate how changes in the reward structure of a task alter visual processing. We use this model to show in principle how experimentally observed attention-dependent changes to visual neuron responses can be interpreted functionally, as a consequence of optimal adaptation toward a given task. In the following sections, we describe the presented stimuli and task, the agent's internal model, and the neural code. Supplementary section 1 (available online) describes the model assumptions in detail and how they influence our results.

**3.1 Visual Detection Task.** In many experimental investigations of goal-directed visual attention, a monkey is instructed (often using a visual cue) that a particular spatial location is task relevant and thus should be attended. In order to receive a reward in the task, the animal is required to make responses that are contingent on stimuli presented at this location, while ignoring distractor stimuli presented at other locations (Luck et al., 1997; Reynolds et al., 2000; Williford & Maunsell, 2006). To capture the main aspects of these experiments, we simulated a visual detection task, in which an agent is presented with one or more stimuli at various locations and has to report whether a stimulus is present at a single target location (see Figure 1). The agent receives a unitary reward for a correct response in the task and no reward otherwise.

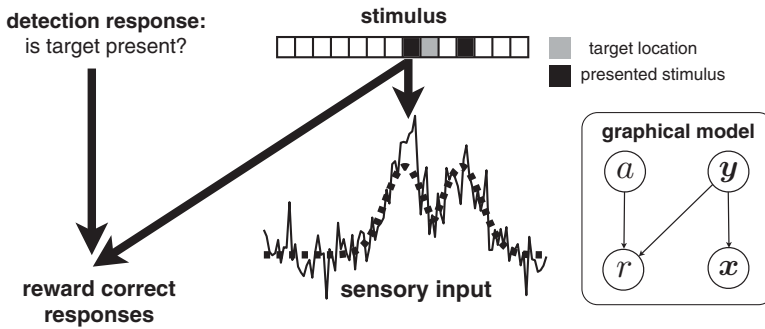


Figure 1: Schematic of the detection task. Presented stimuli ( $y$ ) are represented by binary variables, each indicating whether a stimulus is present at a particular location. Stimuli combine to produce the noisy sensory input ( $x$ ; solid curve). One or more locations are selected as target locations in the detection task (the target is unknown to the agent at the start of the task). The agent gives a response ( $a$ ) indicating whether a stimulus is present at a target location, based on its sensory input and learned model of reward. Correct responses are followed by a reward ( $r$ ). The inset is the corresponding graphical model of the task.

In each attentional condition, stimuli are equally likely to be presented at all locations. The only thing that distinguishes stimuli presented at different locations is whether a reward is delivered for making a detection response. The agent must use this feedback on its performed actions to learn the target location (by adapting the reward model) and to direct attention toward the target (by adapting its sensory model).

The sensory input statistics are described by a binary latent variable model (Puertas, Bornschein, & Lücke, 2010). Presented stimuli are represented by binary hidden variables ( $y_i \in \{0, 1\}$ ), with each variable representing a different spatial location (e.g.,  $y_i = 1$  would indicate that a stimulus is present at the  $i$ th spatial location). There is equal probability for stimuli to be presented at all locations, and stimuli are presented at different locations independent of each other:

$$p(\mathbf{y}) = \prod_{i=1}^{n_y} p(y_i), \quad p(y_i = 1) = \alpha, \quad (3.1)$$

where  $n_y = 20$  denotes the number of spatial locations and  $\alpha$  denotes the probability that a stimulus is presented at any particular location. In our simulations,  $\alpha$  was much smaller than 1, meaning that the visual input statistics were sparse. This sparsity prior was chosen to reflect the statistics of natural images, which are well accounted for by sparse models

(Olshausen & Field, 1996; Berkes, Turner, & Sahani, 2008; see supplementary section 1).

Stimuli combine nonlinearly to generate the sensory input signal received by the agent ( $x$ ), according to

$$x_i = \max_j \{A_{ij}y_j\} + \gamma_i, \quad (3.2)$$

where  $\gamma_i$  is a gaussian noise variable (with zero mean and variance  $\sigma^2 = 0.6$ ) and  $A$  is an  $n_x \times n_y$  matrix of basis functions (we set  $n_x = 20$ ; see supplementary section 1 for a discussion of the nonlinear combination rule).

The basis functions were set up so that a stimulus presented at a single location activates several neighboring inputs. Sensory inputs (components of  $x$ ) were labeled with  $n_x$  equally spaced values between  $-\pi$  and  $\pi$  (producing a vector of spatial locations;  $\tilde{x}$ ). Each of the  $y$ -units (components of  $y$ ) was labeled with  $n_y$  equally spaced values between  $-\pi$  and  $\pi$  (producing a vector of preferred spatial locations; ' $\tilde{y}$ '). Elements of  $A$  were given by

$$A_{ij} = \exp\left(\frac{-(\tilde{x}_i - \tilde{y}_j + 2\pi k)^2}{2\lambda_A^2}\right), \quad (3.3)$$

where  $k$  is an integer, set so that  $-\pi < (\tilde{x}_i - \tilde{y}_j + 2\pi k) < \pi$  (so that the stimulus space is circular and there are no edge effects), and  $\lambda_A$  determines the width of the basis functions (we set  $\lambda_A = 0.35$  for the initial simulations). Columns of  $A$  are plotted in Figure 2. Each plot can be interpreted as the mean sensory activation produced by a stimulus presented at one particular spatial location.

One or more spatial locations, indexed by  $I$ , were chosen as target locations in the task. The detection target ( $t \in \{0, 1\}$ ) was classified as present if a stimulus was present at at least one of the target locations ( $t = 1$  if  $\exists i \in I : y_i = 1$ ). The agent was required to give a response indicating whether it believed the target stimulus was present ( $a = 0$  or  $1$  for a rejection or detection response, respectively). It received a unitary reward for a correct response and no reward otherwise:

$$r = \begin{cases} 0 & \text{if } a \neq t \\ 1 & \text{if } a = t \end{cases} \quad (3.4)$$

**3.2 Agent's Internal Model of Sensory Input.** We assume that the agent uses a hierarchical internal model to infer the hidden causes of the received sensory input (see Figure 3a). Thus, in contrast to the simulated experiment, where spatially localized stimulus features are presented independent of

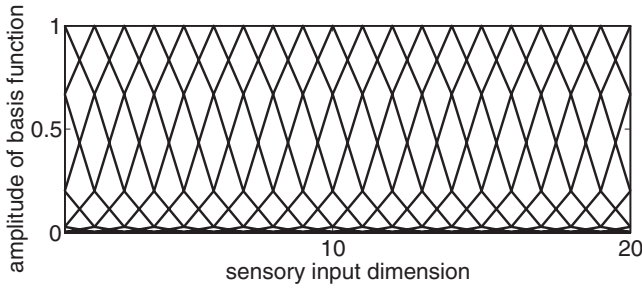


Figure 2: Basis functions used to generate the received sensory input (for the initial simulations, where stimuli included a spatial but not a featural dimension). Each plot shows a single column of the basis function,  $A$ . Individual plots represent the mean sensory input generated by a single active  $y$ -unit. Note that the basis used to generate the sensory input is the same as the agent’s internal model.

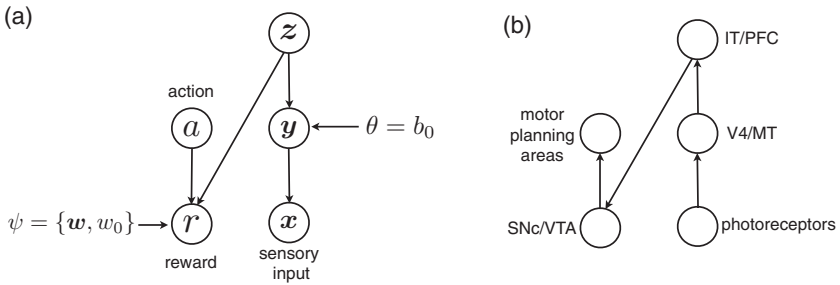


Figure 3: Agent’s internal model of the sensory input and reward. (a) The agent learns a hierarchical model, where high-level hidden variables ( $z$ -units), corresponding to complex spatially distributed image features (e.g., objects or faces), are assumed to determine the state of lower-level hidden variables ( $y$ -units), corresponding to simple spatially localized image features (e.g., orientation or motion direction), which generate the received sensory input ( $x$ ). High-level hidden  $z$ -variables are also assumed to generate the reward received ( $r$ ) for performing different actions ( $a$ ) in the task. During task performance, the agent updates parameters that predict how the reward depends on the high-level hidden variables in its model ( $\psi = \{w, w_0\}$ ), as well as parameters that determine the probability that individual  $y$ -units are active ( $\theta = b_0$ ). (b) Putative mapping of probabilistic model onto neural architecture. Arrows denote the direction of feedforward processing (both direct and indirect). Incoming sensory signals are first processed in low and intermediate visual areas, such as V4 and MT, before being sent to higher-level sensory areas, such as the inferotemporal and prefrontal cortex (IT and PFC). These high-level sensory areas project to regions in the basal ganglia, such as the substantia nigra colliculus (SNc) and ventral tegmental area (VTA), which compute the expected reward for performing different actions.



each other, the agent assumes a higher level of statistical structure, such that certain image features are more likely to be presented together than others.

In the agent's internal model, high-level hidden variables ( $\mathbf{z}$ ) are assumed to generate lower-level hidden variables ( $\mathbf{y}$ ), which give rise to the sensory input ( $x$ ). The joint probability distribution for this model is of the form

$$p(x, \mathbf{y}, \mathbf{z}|\theta) = p(x|\mathbf{y}, \theta)p(\mathbf{y}|\mathbf{z}, \theta)p(\mathbf{z}|\theta), \quad (3.5)$$

where  $\theta$  denotes the parameters of the agent's internal model of the sensory inputs.

All hidden variables are binary ( $y_i \in \{0, 1\}$ ,  $z_i \in \{0, 1\}$ ), while the observed data ( $x$ ) are continuous. For mathematical simplicity, we apply the constraint that a maximum of one  $z$ -unit can be active at a time, with equal probability:

$$p(z_i = 1, z_{j \neq i} = 0|\theta) \propto \rho/n_z, \quad p(\mathbf{z} = 0|\theta) = 1 - \rho, \quad (3.6)$$

where  $n_z = 5$  denotes the number of  $z$ -units in the model and  $\rho$  denotes the probability that one of the  $z$ -units is on (we set  $\rho = 0.5$  for the initial simulations). While the constraint that only one  $z$ -unit may be active at a time may appear quite extreme, it is reasonable that the high-level representation is very sparse, as in general, there will be a small prior probability for any given high-level feature to be present in an image.

Given  $\mathbf{z}$ , the  $y$ -units are assumed to be conditionally independent ( $p(\mathbf{y}|\mathbf{z}, \theta) = \prod_{i=1}^{n_y} p(y_i|\mathbf{z}, \theta)$ ), with a probability of being active given by

$$p(y_i = 1|\mathbf{z}, \theta) = \text{sig}(\mathbf{b}_i^T \mathbf{z} - b_{0i}), \quad (3.7)$$

where  $\text{sig}(x) = (1 + \exp(-x))^{-1}$ ,  $\mathbf{b}_i$  is an  $n_z \times 1$  basis vector and  $b_{0i}$  is a scalar bias term (see section 3.5 for initial values of  $b_{0i}$ , prior to attentional optimization).

The basis vectors  $\mathbf{b}_i$  were set up so that when a given  $z$ -unit is active, there is an increased probability that neighboring  $y$ -units will be active. Components of  $\mathbf{z}$  were labeled with  $n_z$  equally spaced values between  $-\pi$  and  $\pi$  ( $\tilde{\mathbf{z}}$ ). Elements of  $\mathbf{b}_i$  were given by

$$b_{ij} = b_{\max} \exp\left(\frac{-(\tilde{y}_i - \tilde{z}_j + 2\pi k)^2}{2\lambda_B^2}\right), \quad (3.8)$$

where  $\lambda_B$  denotes the width of the basis function (we set  $\lambda_B = 2.5$ ) and  $b_{\max}$  determines how strongly  $z$ -units determine whether the  $y$ -units are on (we set  $b_{\max} = 3$  for the initial simulations). Figure 4 plots the conditional probability that each of the  $y$ -units is on for a given active  $z$ -unit.

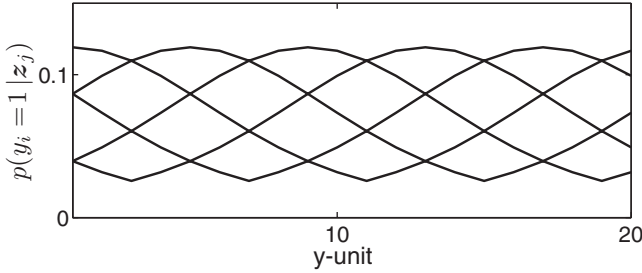


Figure 4: Basis functions used for agents’ internal model of their sensory inputs in the initial simulations, where stimuli included a spatial but not a featural dimension. Each plot shows the probability that the agent assumes different  $y$ -units are active, given a single active  $z$ -unit:  $p(y_i = 1|z_j) = \text{sig}(b_{ij} - b_{i0})$  (before task optimization, with bias terms  $b_{i0}$  set to their initial values).

The agent’s internal model that predicts how the sensory input ( $x$ ) is generated by the hidden causes ( $y$ ) was set identical to the true data generation process described previously (see equations 3.2 and 3.3).

**3.3 Agent’s Internal Model of Reward.** We assume that the agent learns an internal model that predicts how the received reward depends on its performed action and the state of high-level hidden variables in its internal model (see Figure 3). The model of the detection task includes a binary target variable ( $t \in \{0, 1\}$ ), that depends on the state of the  $z$ -units in the agent’s internal model. Given  $z$ , the assumed probability of the target being present is given by

$$p(t = 1|z, \psi) = \text{sig}(\mathbf{w}^T \mathbf{z} - w_0), \tag{3.9}$$

where  $\mathbf{w}$  (an  $n_z \times 1$  vector) and  $w_0$  state how the target variable depends on each of the  $z$ -units, and  $\psi$  denotes collectively the parameters of the agent’s model of reward ( $w_0$  and  $\mathbf{w}$ ). A reward of  $r = 1$  is predicted for a correct response ( $a = t$ ), and no reward ( $r = 0$ ) for an incorrect response ( $a \neq t$ ; see equation 3.4). The agent does not initially know the true location of the detection target:  $\mathbf{w}$  and  $w_0$  have to be learned online through task feedback (see section 3.5).

After receiving a sensory input, the expected reward for reporting that the target is present,  $Q(a = 1; x, \theta, \psi)$ , is equal to the posterior probability that the detection target is present,  $\langle p(t = 1|z, \psi) \rangle_{p(z|x, \theta)}$ . Conversely, the expected reward for reporting that the target is not present,  $Q(a = 0; x, \theta, \psi)$ , is equal to the posterior probability that the target is not present,  $\langle p(t = 0|z, \psi) \rangle_{p(z|x, \theta)}$ .

We assume that the agent makes the response associated with the highest predicted reward. Thus, if the posterior probability that the target is present is greater than 0.5, the agent should make a detection response ( $a = 1$ ); otherwise, the agent should make a rejection response ( $a = 0$ ).

**3.4 Visual Neuron Firing Rates.** Figure 3b illustrates a putative mapping of the probabilistic model used in our simulations onto the neural architecture. The assumed role of the visual system is to infer the posterior probability distribution over the hidden causes. The posterior distribution, encoded in the population activity of visual neurons, is then transmitted to areas of the brain that are responsible for predicting the received reward for performing different actions, allowing the agent to make an appropriate response in the task.

For our simulations, we assume that the mean firing rate of a single visual neuron encodes the posterior probability that a particular hidden cause is active, given the observed sensory input (as in Chikkerur et al., 2010). Thus, the firing rate of the  $i$ th visual neuron is computed directly from Bayes' rule:

$$\begin{aligned} p(y_i = 1|x, \theta) &= \frac{p(x, y_i = 1|\theta)}{p(x|\theta)} \\ &= \frac{\sum_{\mathbf{y}_{/i}} p(x|y_i = 1, \mathbf{y}_{/i}, \theta)p(y_i = 1, \mathbf{y}_{/i}|\theta)}{\sum_{\mathbf{y}} p(x|\mathbf{y}, \theta)p(\mathbf{y}|\theta)}, \end{aligned} \quad (3.10)$$

where  $\mathbf{y}_{/i}$  represents a vector of all the components of  $\mathbf{y}$ , except for the  $i$ th component, and the summation is taken over all possible hidden states.

For our simulations, there were sufficiently few latent variables that we were able to perform the summation over the latent states directly. However, if there is a large number of hidden variables, this summation will become intractable, and an approximate algorithm must be used. Shelton, Bornschein, Sheikh, Berkes, and Lücke (2011) describe a biologically plausible algorithm that could be used to perform approximate inference on a binary latent variable model similar to the one used in our simulations (Puertas et al., 2010).

The stimulus selectivity of a given neuron is largely determined by the basis function of the hidden variable it encodes. In other words, if the hidden variable encoded by a given neuron typically generates a specific profile of sensory activity, then receiving this same sensory activation profile will imply that the hidden cause is active and the neuron will respond with a high firing rate. The basis functions used in our simulations were spatially localized, so that model neurons responded most strongly to stimuli presented at a small number of neighboring locations (their receptive field, RF). The basis functions of the low-level  $y$ -units were narrower than the

basis function of the high-level  $z$ -units (compare Figures 2 and 4), so that neurons encoding  $y$ -units had smaller RFs than neurons encoding  $z$ -units. Note, however, that in general, a neuron's RF is not identical to the basis function of the encoded variable. Although basis functions are an invariant property of the generative model, the measured RF will depend on the types of stimuli presented.

**3.5 Task Optimization.** We hypothesized that attentional processes continually adapt the agent's internal model to improve their predictions of the received reward (at the potential cost of learning a worse internal model of the received sensory inputs). Parameters of the agent's internal model ( $\theta$  and  $\psi$ ) are adapted online to maximize the log probability of the received reward. After each trial, model parameters are updated according to

$$\theta_{new} \leftarrow \theta + \eta_i \partial_\theta l_i(\theta, \psi), \quad \psi_{new} \leftarrow \psi + \eta_i \partial_\psi l_i(\theta, \psi), \quad (3.11)$$

where  $l(\theta, \psi) \equiv \log p(r|x, x, \theta, \psi)$ , and  $\eta$  is the learning rate. In supplementary section 2, we show that the derivative of the online objective function can be written as

$$\partial_\psi l_i(\theta, \psi) = \langle \partial_\psi \log p(r_i|a_i, \mathbf{z}, \psi) \rangle_{p(\mathbf{z}|x_i, r_i, a_i, \theta, \psi)}, \quad (3.12)$$

$$\begin{aligned} \partial_\theta l_i(\theta, \psi) &= \langle \partial_\theta \log p(\mathbf{y}, \mathbf{z}, x_i|\theta) \rangle_{p(\mathbf{y}, \mathbf{z}|x_i, r_i, a_i, \theta, \psi)} \\ &\quad - \langle \partial_\theta \log p(\mathbf{y}, \mathbf{z}, x_i|\theta) \rangle_{p(\mathbf{y}, \mathbf{z}|x_i, \theta)}. \end{aligned} \quad (3.13)$$

For parameters to converge on stable values, we used a learning rate that decreased as a function of the trial number, according to  $\eta = \eta_0/(1 + i/n_0)$  (where  $i$  is the trial number and  $\eta_0$  and  $n_0$  are parameters that determine the initial learning rate and how fast it decays, set to 0.05 and  $10^4$ , respectively). Learning was terminated after  $10^5$  trials, when the model parameters were observed to converge on stable values.

We postulated that over the short timescales associated with visual attention, only the prior probability that individual hidden  $y$ -units are active varies (determined by the bias terms,  $b_{0i}$ , in equation 3.7), while other aspects of the internal model are unchanged. The gradient of the objective function used to update  $b_{0i}$  is given by

$$\langle \partial_{b_{0i}} \log p(x, \mathbf{y}, \mathbf{z}) \rangle = \langle \text{sig}(\mathbf{b}_i^T \mathbf{z} - b_{0i}) - y_i \rangle. \quad (3.14)$$

Note that evaluating this expression requires computing only first-order statistics, such as the mean activation of the  $y$ -units. In comparison, updating the basis functions ( $\mathbf{b}_i$  and  $\mathbf{A}$ ) would require computing second-order statistics, which are harder to estimate from a limited supply of noisy data.

We initialized the bias term terms,  $b_{0i}$ , to take equal values, such that the prior probability that each  $y$ -unit was active was exactly equal to the true probability that a stimulus was presented at each location ( $\alpha = p(y_i = 1 | \theta_{init})$ ). Consequently, before optimization, the only difference between the agent's internal model and the true model describing how the sensory inputs were generated was related to the second-order statistics describing the probability that stimuli were presented at different locations at the same time. For the true model, all  $y$ -units were independent, while for the agent's internal model, there was a higher probability that adjacent  $y$ -units were simultaneously active.

While the agent is assumed to know the general structure of the task (i.e., that they receive a unitary reward for detecting a visual target), they do not know in advance where the target is (both  $w_0$  and  $\mathbf{w}$  are set to zero initially). These parameters ( $\psi \equiv \{w_0, \mathbf{w}\}$ ) are learned online on the basis of the reward received for performing different actions. The objective function gradient used to update  $\mathbf{w}$  and  $w_0$  is given by

$$\langle \partial_{w_i} \log p(r|a, z) \rangle = \langle z_i (r - \text{sig}(\mathbf{w}^T \mathbf{z} - w_0)) \rangle, \quad (3.15)$$

$$\langle \partial_{w_0} \log p(r|a, z) \rangle = -\langle r - \text{sig}(\mathbf{w}^T \mathbf{z} - w_0) \rangle. \quad (3.16)$$

At the beginning of each task, we initialized  $w_0 = 0$  and  $\mathbf{w} = 0$ , implying that the agent had no knowledge of the location of the detection target.

Note that our aim was to investigate the effects of attentional optimization rather than the temporal dynamics of the optimization process itself. Thus, while we assume that attentional modulation of visual neuron responses is learned online from task feedback, in reality, the attentional state could also be altered more quickly, based on information received from visual cues or previous experience in the task (see section 5).

## 4 Results

---

**4.1 Attentional Modulation of Detection Performance.** We first asked how attention altered performance in the detection task. We considered two conditions: a no-attention condition, in which the agent optimized its reward model but not its sensory model, and an attend-target condition, where the agent optimized both its reward model and its sensory model.

On each trial, the agent estimated the probability that a stimulus was present at a target location,  $p(t = 1|x)$ , to decide whether to make a detection response. We used the agent's estimates of  $p(t = 1|x)$  to plot receiver operating characteristic curves (ROC) for each attentional condition (see Figure 5a). The area under the ROC curve (the AUC) provides a measure of detection performance that is independent of the threshold used for classification: an AUC value of 1 indicates perfect performance, while an AUC value of 0.5 indicates chance performance (Fawcett, 2006). As expected,

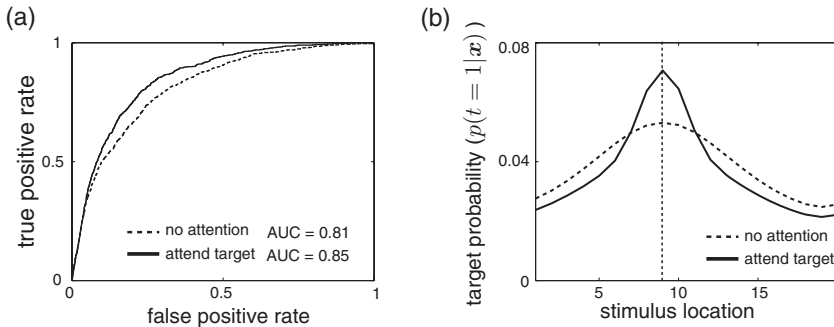


Figure 5: Attention-dependent changes to detection performance. (a) Receiver operating characteristic curves, indicating how well the agent is able to classify whether the detection target is present. The area under each curve (AUC values) gives a summary statistic of how well the agent is able to classify whether the detection target is present. The agent does marginally better in the attend-target condition than in the no-attention condition. (b) The estimated probability that a stimulus is present at the target location ( $p(t = 1|x)$ ), plotted as a function of presented stimulus location. In the attend-target condition, the agent is better able to discriminate whether a stimulus is presented at the target location.

detection performance was better in the attend-target condition (AUC = 0.85) than in the no-attention condition (AUC = 0.81). The magnitude of this performance increase was observed to be highly dependent on the precise setup of the task (e.g., increasing the sensory noise leads to larger attention-dependent improvements in performance). However, while the magnitude of attention-dependent changes to performance varied depending on the task, the qualitative effect of attention was always the same: to improve performance in the detection task.

To understand how attention alters detection performance, we plotted the estimated probability that a stimulus was present at a target location ( $p(t = 1|x)$ ) versus the true stimulus location (see Figure 5b). In the attend-target condition, the agent's estimates of  $p(t = 1|x)$  were increased for stimuli close to the target location and reduced for stimuli far from the target location. Thus, in the attend-target condition, the agent was better able to detect stimuli at the target location, while ignoring stimuli at other locations.<sup>1</sup>

**4.2 Attentional Modulation of Neural Population Response.** We next asked how attention alters the internal sensory representation, encoded by

<sup>1</sup>Note that variations in the agent's estimates of  $p(t = 1|x)$  matter more than its baseline value, as changes in baseline can be easily compensated for by varying the detection threshold.

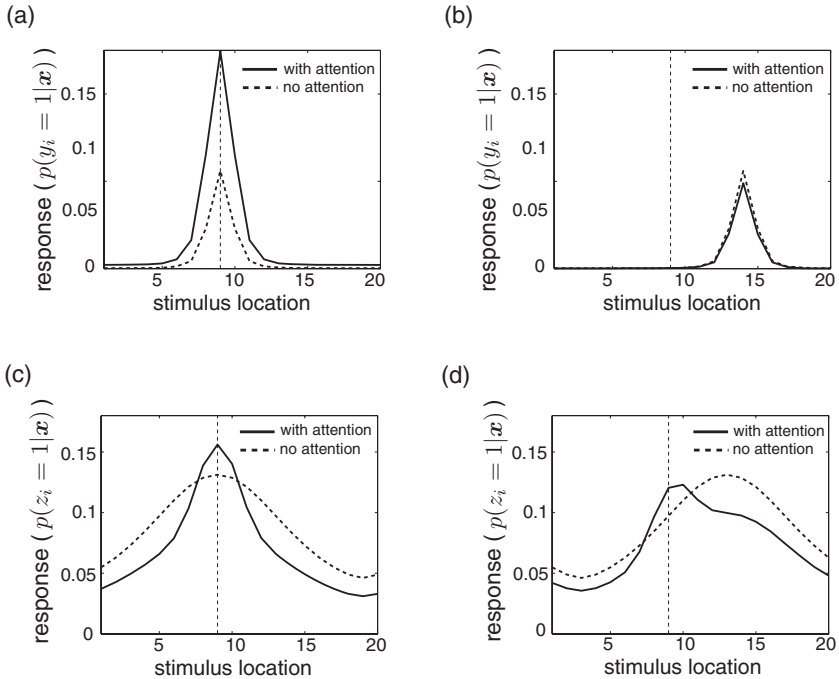


Figure 6: Influence of spatial attention on the spatial tuning curves of midlevel and high-level visual neurons. (a) Spatial tuning curve of a midlevel neuron tuned to the target location (indicated by a vertical dashed line). The neuron’s response is plotted as a function of the presented stimulus location, with or without attention directed toward the target location. (b) Spatial tuning curve of a midlevel neuron tuned elsewhere. (c) Spatial tuning curve of a high-level neuron tuned to the target location. (d) Spatial tuning curve of a midlevel neuron tuned elsewhere.

the visual neuron responses. The model was set up so that midlevel neurons (encoding  $y$ -units in the agent’s internal model) were highly sensitive to the presented stimulus location, with each neuron responding only to stimuli presented near the neuron’s preferred location (see Figures 6a and 6b, dashed line). In contrast, high-level neurons (encoding  $z$ -units in the agent’s internal model) were relatively insensitive to the presented stimulus location (see Figures 6c and 6d, dashed line).

In our model, the agent relied on the responses of high-level neurons to choose which action to perform. However, as high-level neurons were insensitive to the stimulus location, the agent was not able to discriminate between stimuli presented at task-relevant and task-irrelevant locations, impairing its performance in the task. Following attentional optimization

toward the task (see section 3.5), the agent learned to associate increased prior probability for stimuli at the target location. This learned prior did not reflect the true stimulus statistics (stimuli were equally likely at each location) but instead compensated for the mismatch between the agent's internal model and the true structure of the task.

The attentional prior increased the gain of midlevel neurons whose preferred location was near the target location (see Figure 6a) while decreasing the gain of neurons whose preferred location was far from the target location (see Figure 6b). This change in the gain of midlevel neurons resulted in changes to the stimulus selectivity of high-level neurons, which became differentially more sensitive to stimuli presented at the target location (see Figures 6c and 6d). The net result was that in the attend-target condition, high-level neural responses were a better predictor of whether a stimulus was present at the target location, allowing the agent to improve its task performance.

**4.3 Attentional Modulation of the Contrast Response Function.** There have been a number of controversies about how goal-directed attention alters sensory neural responses. A prominent example is attention-dependent changes to the firing rates of V4 neurons with varying stimulus contrast. Previous experiments have reported very different findings. Williford and Maunsell (2006) observed a “response gain” effect, with increases in neural firing rates for all stimulus contrasts, while Reynolds et al. (2000) observed a “contrast gain” effect, consistent with an increase in the effective stimulus contrast. Reynolds and Heeger (2009) proposed a phenomenological model to account for these differences, proposing that they are due to variations in the relative size of the focus of attention and the stimulus between experiments: a narrow focus of attention would give rise to a response gain effect, while a broad focus of attention would give rise to a contrast gain effect. We use our normative model to ask why attention might alter neural responses in this way.

To manipulate the size of the attentional focus, we varied the number of target locations in the detection task. We simulated two experimental conditions: one with a single target location (narrow attentional focus) and another with multiple neighboring target locations (broad attentional focus, with seven neighboring locations chosen as targets). Note that only the reward contingencies changed for the different attentional conditions; the stimulus statistics were always the same (see section 3.1).

In the narrow attentional focus condition, the agent learned to associate an increased prior probability that hidden causes representing stimuli at this location were active (see Figure 7a). In the broad attentional focus condition, there was a broader change in its learned prior, with increases in the prior probability for hidden causes representing all of the target locations (see Figure 7b).



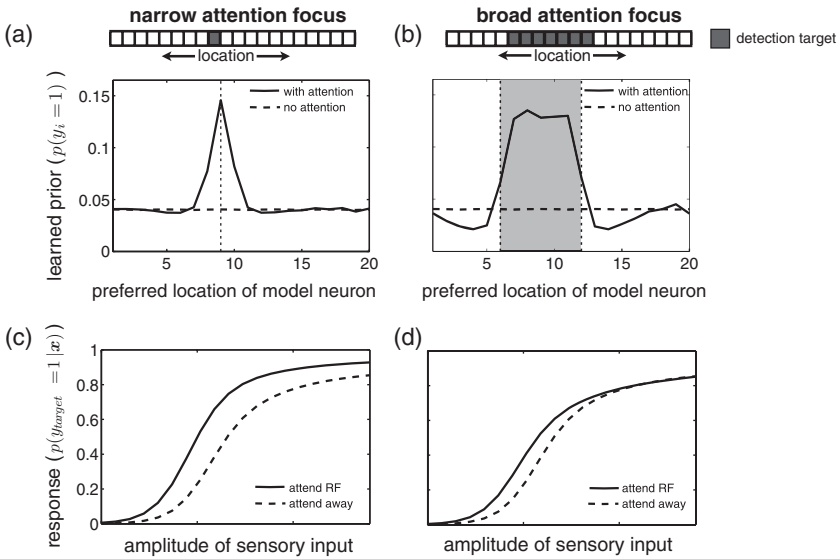


Figure 7: Attentional modulation of neural contrast response function. (a, b) Prior probability assumed by the agent that each of the hidden causes is active, without attention (dashed line) or with either a narrow (a) or a broad (b) focus of attention (solid line). The attended spatial region is represented by the shaded area. (bottom panels) Model neuron response, as a function of the amplitude of a sensory input at the preferred location, without attention (dashed line) or with either a narrow (c) or a broad (d) focus of attention (solid line).

Neurons in visual area V4 were hypothesized to encode information about hidden variables at an intermediate level of the agent’s internal model (i.e., components of  $\mathbf{y}$ ). To obtain neural contrast response functions (CRFs), we plotted the mean firing rate of a model neuron while varying the amplitude of a sensory input centered at its preferred location ( $x = ca_i$ , where  $a_i$  is the  $i$ th column of  $\mathbf{A}$ , and  $c$  represents the stimulus contrast). The resulting CRF was qualitatively similar to experiment, increasing monotonically at intermediate sensory input amplitudes, before saturating at high amplitudes. The effect of spatial attention was consistent with Reynolds and Heeger’s (2009) divisive normalization model: directing a narrow focus of attention toward the presented stimulus location increased the response of a neuron tuned to this location for all sensory input amplitudes; a broad focus of attention increased the response of this neuron only at intermediate sensory input amplitudes (see Figures 7c and 7d, respectively).

**4.4 Comparison with Normalization Model of Attention.** To see why attention alters neural CRFs as it does in our model, consider the expression

for neural firing rates, computed directly from Bayes' law (see equation 3.10). Because the image statistics are sparse, meaning that there is a very small probability that multiple image features are present in any given image, we can approximate the response of the  $i$ th neuron by

$$p(y_i = 1|x) \approx \frac{p(x|y_i)p(y_i)}{p(x|y_0)p(y_0) + \sum_{j=1}^{n_y} p(x|y_j)p(y_j)}, \quad (4.1)$$

where  $y_i$  denotes a hidden state with only one active  $y$ -unit (i.e.,  $y_i \equiv (0, \dots, 0, 1, 0, \dots, 0)$  with only  $y_i = 1$ ), and  $y_0$  denotes a hidden state with all  $y$ -units inactive (i.e.,  $y_0 = 0$ ). This expression shows that divisive normalization of neural firing rates is predicted in our model as a consequence of performing Bayesian inference on a sparse binary latent variable model. Here, divisive normalization comes about due to a well-known Bayesian phenomenon called explaining away, in which different hidden causes compete with each other to explain the observed sensory input.

We can rewrite the expression for the neural firing rates as

$$f_i(x) \sim \frac{\mathcal{A}_i E_i(x)}{1 + \sum_{j=1}^{n_y} \mathcal{A}_j E_j(x)}, \quad (4.2)$$

where  $E_i(x) \propto \exp\left(\frac{a_i^T x}{\sigma^2}\right)$  and  $\mathcal{A}_i = \frac{p(y_i)}{p(y_0)}$ . Attention alters the prior probability that the individual  $y$ -units are on, increasing the value of  $\mathcal{A}_i$  for neurons that are tuned to attended stimuli.  $E_i(x)$  is determined by the sensory input alone and does not depend on the attentional state of the agent.

At low contrasts, neural firing rates can be approximated by  $f_i(x) \sim \mathcal{A}_i E_i(x)$ , so that both a narrow and a broad focus of attention alter neural responses multiplicatively, increasing the firing of neurons that are tuned to attended spatial locations. At high contrasts, neural firing rates can be approximated by  $f_i(x) \sim \frac{\mathcal{A}_i E_i(x)}{\sum_{j=1}^{n_y} \mathcal{A}_j E_j(x)}$ . In this case, a broad focus of attention produces similar increases to both the numerator and the denominator, so that the response of a neuron that is tuned to an attended stimulus is unchanged by attention (see Figure 7d). A narrow focus of attention increases the numerator by a larger factor than the denominator, so that the response of a neuron that is tuned to an attended stimulus is increased (as in Figure 7c).

Both the expression for neural firing rates and the modulatory effects of attention in our model are similar to Reynolds and Heeger's (2009) normalization model of attention. However, while divisive normalization was an ad hoc assumption in Reynolds and Heeger's model, in our work it comes about as a direct consequence of performing Bayesian inference on a particular form of internal model. Likewise, while Reynolds and Heeger

specified an attention field, which multiplicatively scaled the gain of the feedforward excitatory input to the network, in our work, attentional modulation of neural responses comes about as a result of optimization toward the task and is thus entirely determined by the behavioral task and the agent's internal model.

**4.5 Attentional Modulation of Sensory Tuning Curves.** We investigated how goal-directed attention alters neural tuning curves in our model. To do this, we extended our model to include both a featural and a spatial dimension. We altered the basis functions that determined the image features represented by the hidden units, so that each model neuron (corresponding to a component of  $\mathbf{y}$ ) was selective to both a stimulus feature (e.g., orientation, or motion direction) and a spatial location.

Every sensory input (component of  $\mathbf{x}$ ) was allocated a feature label (consisting of two lists of  $n_x/2$  equally spaced values between  $-\pi$  and  $\pi$ ;  $\tilde{\mathbf{x}}_1$ ), and a spatial label ( $n_x/2$  lists of two spatial locations,  $(0, \pi)$ ;  $\tilde{\mathbf{x}}_2$ ). The  $\mathbf{y}$ -units were labeled in the same way: each with a corresponding spatial location and feature ( $\tilde{\mathbf{y}}_1$  and  $\tilde{\mathbf{y}}_2$ , respectively). Elements of  $\mathbf{A}$  were given by

$$A_{ij} = \exp \left( \frac{-(\tilde{x}_{i1} - \tilde{y}_{j1} + 2\pi k)^2}{2\lambda_{ftr}^2} + \frac{-(\tilde{x}_{i2} - \tilde{y}_{j2} + 2\pi k)^2}{2\lambda_{spt}^2} \right), \quad (4.3)$$

where  $\lambda_{ftr}$  and  $\lambda_{spt}$  are parameters determining the width of the basis function along feature and spatial dimensions, respectively. The basis functions for the z-units were calculated in the same way as for the previous simulations (see section 3.2), with all z-units allocated a feature but not a spatial label (i.e.,  $\tilde{\mathbf{y}}_i$  was replaced by the feature label,  $\tilde{y}_{i1}$ ). We set,  $\lambda_{ftr} = 1.2$  and  $\lambda_{spt} = 2$  (see the next section for a discussion of how we set  $\lambda_{spt}$ ). We also increased the model sparsity, setting  $\rho = 0.3$  and  $b_{\max} = 2$  (so that  $p(y_i = 1 | \theta_{init}) \approx 0.02$ ). This increase in sparsity was required to produce robust surround suppression for the simulations described in the next section (but was not critical for simulating the feature tuning curves).

We simulated two experimental conditions. In the first condition (spatial attention), one of two spatial locations was selected as a target in the detection task. In the second condition (feature-based attention), only certain features were chosen as targets. Spatial attention caused the agent to associate a high prior probability that hidden variables representing the attended location were active, but a uniform prior probability that hidden variables representing different features were active (see Figure 8a). Conversely, feature-based attention caused the agent to associate a high prior probability that hidden variables representing attended features were active, but a uniform prior probability that hidden variables representing both spatial locations were active (see Figure 8b).

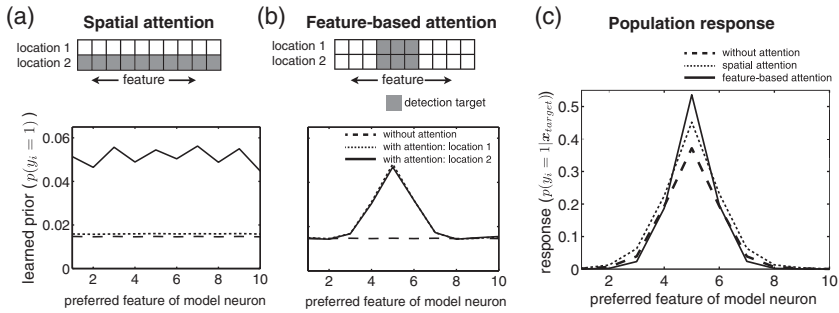


Figure 8: Influence of spatial and feature-based attention on the population response. (a, b) Prior probability assumed by the agent that each of the hidden causes is active, without attention (solid line) or with spatial (a) or feature-based (b) attention. (c) Neural population response in the absence of attention (dashed line) or with attention directed toward the presented stimulus feature (solid line) or spatial location (dotted line).

Attending to the presented stimulus location increased the responses of neurons tuned to this location, with no sharpening in the population response (see Figure 8c, dotted line). Similar effects have been observed experimentally in visual area V4 when attention is directed toward a particular spatial location (McAdams & Maunsell, 1999). In contrast, we found that attending to the presented stimulus feature produced a sharpening in the population response; the responses of model neurons that were selective for the attended feature were most strongly increased by attention (see Figure 8c, solid line). Martinez-Trujillo and Treue (2004) reported a similar effect in visual area MT when animals were directed toward a particular motion direction. Our results are also consistent with Reynolds and Heeger's (2009) normalization model of attention.

Also consistent with the experimental findings of Martinez-Trujillo and Treue (2004), our model predicted a small suppression in the responses of model neurons tuned to unattended features. In our model, this suppression came about because the agent accorded greater probability to the possibility that the sensory input was produced by hidden causes representing attended features, at the expense of a reduction in the probability that it was produced by hidden causes representing other, unattended, features.

Experimentally it has been shown that attention-dependent suppression of neural responses is particularly strong when there are multiple stimuli within the cell's RF (Moran & Desimone, 1985; Reynolds et al., 1999). Although we do not explicitly model this effect, it is easy to see how it could come about for our model. When there is one stimulus within a cell's RF, directing attention away from or toward the presented stimulus will

induce a multiplicative change to the neuron's response by altering the numerator in equation 4.2. When two stimuli are present within the cell's RF, attending toward one of the stimuli will also alter suppression that comes from the other stimulus via the denominator in equation 4.2, resulting in larger changes in the neuron's response. This effect was demonstrated by Reynolds and Heeger (2009) in their normalization model of attention.

**4.6 Attentional Modulation of Center-Surround Interactions.** The responses of neurons in the visual cortex are modulated by stimuli located outside their classical RF that do not evoke a response when presented alone. Typically, presenting a stimulus outside a neuron's RF suppresses its response, compared to when there is only a single stimulus presented within its RF, a phenomenon called surround suppression (Seriès, Lorenceau, & Frégnac, 2003). Sundberg et al. (2009) found that in visual area V4, attending to a stimulus located within the RF reduces the suppressive influence of a stimulus presented at the surround, while attending to the surround increases this suppression.

We used the setup described in the previous section to measure the degree of surround suppression in our model in the absence of attention or with attention directed to either the RF center or the surround (see Figure 9a). By definition, a stimulus in the RF surround should not elicit a response when presented alone, although it may suppress the response of a neuron to a stimulus simultaneously presented in the RF. To reproduce this behavior in our model, we needed to specify the spatial width of the basis functions, determined by  $\lambda_{spt}$  in equation 4.3. If they are too broad, surround stimuli elicit a response when presented alone; too small, and there is no surround suppression (we found that  $\lambda_{sptl} = 2$ , produced the required behavior; see Figure 9a).

Directing attention toward the RF increased the model neuron response toward a single stimulus presented within the RF, while decreasing the suppression from a second stimulus presented at the surround (see Figures 9b and 9c). Directing attention to the surround did not significantly alter the model neuron response when a single stimulus was presented within the RF, but did increase the suppression caused by a second stimulus presented at the surround (see Figure 9b and 9c, left panel). In both conditions, the response of the model neuron to a stimulus presented at the surround alone was negligible. Qualitatively similar results were obtained by Sundberg et al. (2009) (see Figure 9c, right panel).

Note that in all attentional conditions, surround suppression was significantly stronger in our model than in the population-averaged data (by a factor of 2). However, the important qualitative aspect of the data that we sought to capture was the effect of attention on surround suppression rather than the absolute magnitude of surround suppression. Indeed, while the

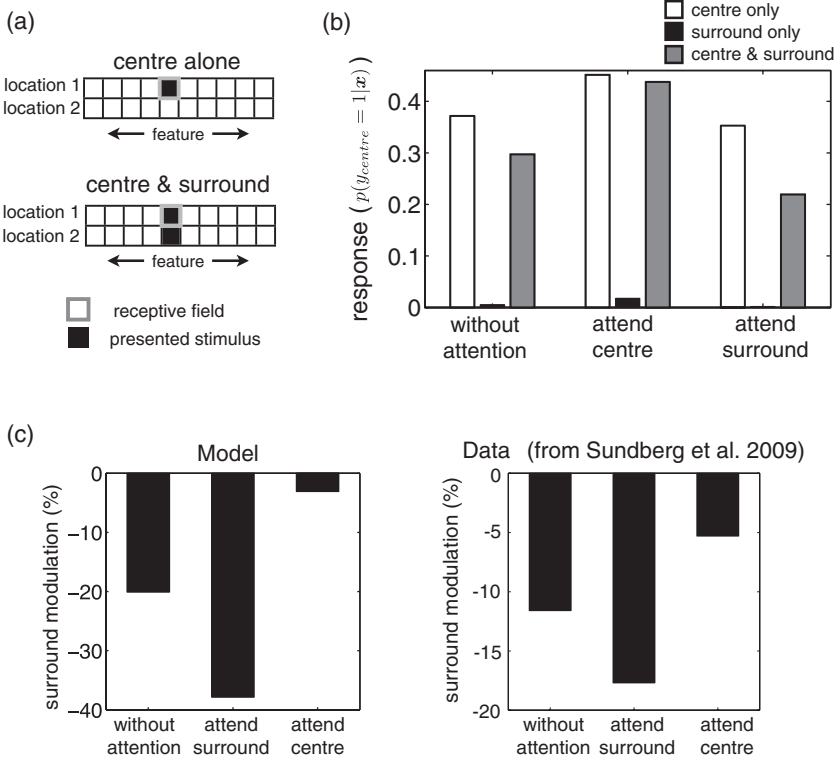


Figure 9: Attentional modulation of center-surround suppression. (a) Schematic of test stimuli. Neural responses were measured with either a single stimulus presented at their RF (top) or with stimuli presented at both their RF center and surround (bottom). (b) Response of a model neuron to a stimulus presented in the RF center (white), surround (black), or both the RF center and surround (gray), without attention, or with attention directed toward the RF center (attend centre) or surround (attend surround). (c) Fractional change in model neuron response when a second stimulus is presented in the surround for each of the three attentional conditions. The left panel shows the predictions of our model, while the right panel shows the population-averaged data from area V4, in an experiment conducted by Sundberg et al. (2009).

qualitative effects of attention were robust to changes in model parameters, the absolute magnitude of surround suppression depended on our choice of model parameters (and experimentally, Sundberg et al., 2009, observed a large variability in the degree of surround suppression across different neurons).

## 5 Discussion

---

We extended previous Bayesian models of visual processing (Hyvärinen, 2010) to account for the effects of behavioral demands on visual neuron responses, hypothesizing that the brain learns a probabilistic model that predicts how both the sensory input and reward received for performing different actions are determined by a common set of hidden causes (Sahani, 2004). We developed a simple model of visual processing to show in principle how our proposed framework can be used to make concrete predictions about how task-dependent attention modulates visual neuron responses. Our framework has two main advantages. First, it has predictive power: in theory, changes to neural responses can be predicted as a direct consequence of the presented stimuli and behavioral task. Second, predicted changes to neural responses have a direct functional meaning: they correspond to changes in the believed causes of the sensory input.

In order to make concrete predictions about the effects of attention on visual neuron responses, we needed to make certain assumptions about the agent's internal model of its environment. First, we assumed that the agent learns a model in which binary hidden causes are responsible for generating its received input. This internal model was very similar to a previous work by Puertas et al. (2010). Second, we assumed that the agent's internal model was sparse, meaning that there was a small prior probability for any particular hidden cause to be active (Olshausen & Field, 1996, 1997). This sparsity prior leads to competition between different possible causes of the sensory input, and in our neural model results in surround-suppression of neural responses. Third, we assumed that the agent performs inference on a hierarchy of image features (Karklin & Lewicki, 2003), and that its behavioral responses depended on only high-level hidden variables in their internal model. In our neural model, this corresponds to relying on the responses of high-level neurons with large receptive fields to choose which action to perform. Finally, we assumed that attention alters the bias terms in the agent's internal model but not the basis functions. This corresponds to altering the gain of individual neurons but not the network connectivity (Dayan & Zemel, 1999; Yu, Dayan, & Cohen, 2009). Most of our assumptions are not new but correspond to assumptions made implicitly in many phenomenological and mechanistic models of attention (Reynolds & Heeger, 2009; Ghose, 2009; Lee & Maunsell, 2009). However, in contrast to these models, we justify our assumptions from functional principles to provide insight into why attention alters visual neuron responses as it does.

The predictions and mathematical formulation of our model bear strong similarities to the normalization model of attention, proposed by Reynolds and Heeger (2009). Recently, both Schwartz and Coen-Cagli (2013) and Chikkerur et al. (2010) showed that Reynolds and Heeger's normalization model can be derived using a Bayesian framework. In their models, attention is hypothesized to modulate the agent's perceptual prior

(Chikkerur et al., 2010) or the feedforward inputs to neurons at attended locations (Schwartz & Coen-Cagli, 2013). However, in both models, attention-dependent changes are specified explicitly, without stating why these changes might come about. As a result, these models suffer from the same limitation as Reynolds and Heeger's normalization model: they do not explain how attention should be shaped by behavioral demands and sensory experience. In contrast, in our model, task-dependent changes to the agent's internal model are learned automatically by the agent in order to improve their predictions of the received reward.

Several studies have tried to explain visual attention in normative terms, under the hypothesis that it corresponds to changes in the perceptual prior (Dayan & Zemel, 1999; Rao, 2005; Chikkerur et al., 2010; Yu & Dayan, 2005; Yu et al., 2009). However, in the absence of any changes to the presented stimulus statistics, it is not clear why the perceptual prior should be altered by task demands. Indeed, in nearly all Bayesian models of attention, changes to the agent's prior are either specified explicitly (Dayan & Zemel, 1999; Rao, 2005; Chikkerur et al., 2010; Whiteley & Sahani, 2012), or learned directly from the stimulus statistics (Yu & Dayan, 2005; Yu et al., 2009). Here, we show that in certain circumstances, it is desirable to alter the perceptual prior, even in the absence of any changes to the stimulus statistics. In our proposed framework, the agent continuously adapts the internal model to improve predictions of the reward associated with each action. When there is a mismatch between the agent's internal model and the true structure of the task, improvements in their predictions of reward may come at the expense of learning a worse model of the sensory input statistics. Consequently, their learned prior will differ from the true stimulus statistics (as in Figures 7a and 7b).

In our simulations, we implemented a specific type of model mismatch in which the stimulus features relevant to the task (the detection targets) are smaller than the features used to decide which action to perform. This is analogous to previous modeling work, in which the agent uses the response of neurons with large RFs to detect stimuli presented in a small task-relevant region of space (Yu et al., 2009; Dayan & Solomon, 2010; Dayan & Daw, 2008; Liu, Yu, & Holmes, 2009). In this work, perceptual performance is limited because neurons integrate sensory signals from both task-relevant and task-irrelevant spatial locations. Attention improves performance by selectively boosting neural inputs that are selective to stimuli at task-relevant locations. Previous authors suggested that perceptual performance is constrained in this way because of the limited number (and thus, necessarily large size) of neural RFs available to cover the visual scene (Dayan & Daw, 2008; Dayan & Zemel, 1999). However, this cannot explain why the agent does not use information encoded by low-level visual neurons with small RFs to perform the task. Here, we propose an alternative explanation: that perceptual performance is constrained by the need to learn a simple behavioral strategy that can be quickly altered in response to changing behavioral demands.



One way to achieve this goal could be to learn a simple mapping between the responses of a small number of high-level neurons (with large RFs) and the reward associated with each action.

Recently, Whiteley and Sahani (2012) proposed that in complex environments, the agent simplifies perceptual inference by using an approximate internal model that neglects statistical dependencies between stimuli. This results in a mismatch between the agent's internal model and its external environment, which reduces its perceptual performance. Whiteley et al. hypothesized that attention compensates for this reduction in performance, forming part of an approximate inference algorithm that selectively improves perceptual accuracy for certain attended features or stimuli.

In both our model and the model of Whiteley and Sahani (2012), attention is required because of a mismatch between the agent's internal model and its environment. In Whiteley's model, this mismatch occurs because the agent neglects dependencies between hidden variables; in our model, it occurs because the agent learns a simplified model of the task, in which the received reward is assumed to depend on a limited number of high-level variables. Experimentally one could distinguish between these different scenarios by investigating when attention is most strongly recruited: when there are complex statistical dependencies between stimuli (as predicted by Whiteley and Sahani's model) or when task-relevant stimulus features are localized in a particular spatial or featural dimension (as predicted by our model). However, rather than there only ever being one type of model mismatch, it is more likely that attention is required in a range of different situations to compensate for different mismatches between the agent's internal model and its external environment. Put in this broader context, we believe that Whiteley and Sahani's model is not incompatible with our framework. For example, one could imagine a hybrid of both models in which reward feedback is used to determine which stimulus features are task relevant, controlling an approximate inference algorithm that improves perceptual accuracy toward these features.

In this letter, we focused on attentional modulation of midlevel neural responses. However, our model also predicts how attention should modulate the responses of higher-level visual neurons. In our simulations, attention dynamically alters the RF profiles of high-level neurons, shrinking them around attended stimuli (see Figure 6c) or shifting their centers toward attended locations (see Figure 6d). This prediction is supported by experimental recordings in area MT, which observe dynamic reshaping of neural RFs as a result of visual attention (Womelsdorf et al., 2006). Of course, in the brain, there is no clear demarcation between high-level or midlevel neurons. However, in the context of our model, what matters is the ratio between a neuron's RF size and the size of task-relevant stimulus features. A neuron is considered to be "high level" if its RF is significantly larger than the task-relevant stimulus features. (Note that while we discuss only

spatial attention here, an analogous argument could be made in the feature domain to describe feature-based attention.)

At the behavioral level, our model predicts that stimuli should be perceived as being more similar to attended stimuli than they actually are. This is because attention-dependent changes to the perceptual prior will induce an estimation bias toward task-relevant stimulus features. While estimation biases have been observed experimentally in response to changes in the presented stimulus statistics (Chalk, Seitz, & Serié, 2010), we predict that they should also be induced by changes to the behavioral task alone. Experimentally, different behavioral tasks have been found to give rise to qualitatively different types of perceptual bias. For example, Jazayeri (2007) found that after performing a discrimination task with visual motion stimuli, subjects report stimuli as moving farther away from the discrimination boundary than they actually are. Therefore, while we simulated a simple detection task, it would be interesting to use our modeling framework in the future to investigate how different behavioral tasks and stimuli influence perception.

At present, it is unknown how probability distributions are represented in the brain (Fiser, Berkes, Orbán, & Lengyel, 2010; Shelton et al., 2011; Deneve, 2008; Ma, Beck, Latham, & Pouget, 2006). A current area of debate is whether neural firing rates encode samples from a probability distribution (Fiser et al., 2010; Shelton et al., 2011); or parameters, such as the mean and variance of the distribution (Deneve, 2008; Ma et al., 2006). In our simulations, we assumed that mean firing rates are proportional to the probability that individual hidden causes contributed to generating the received sensory input. While this coding scheme was chosen for simplicity, it produces mean firing rates that are qualitatively consistent with a sampling code (Shelton et al., 2011). Meanwhile, certain parametric codes, such as the coding scheme proposed by Deneve (2008), predict mean firing rates that are qualitatively consistent with our model (i.e., they scale monotonically with the posterior probability that encoded latent variables are active).

We investigated short-term effects of behavioral context, focusing specifically on visual attention. We hypothesized that over these timescales, only the sensitivity of individual neurons (the prior) varies, while the network connectivity (the basis functions) remains constant. This restriction could be removed to investigate changes that take place over longer timescales. Currently, the relationship between different types of sensory learning—for example, attentional (Eckstein, Abbey, Pham, & Shimozaki, 2004; Jiang & Chun, 2001) versus perceptual learning (Fahle, 2005; Seitz, Kim, & Watanabe, 2009)—and how they depend on the training paradigm, is an active area of research. Our framework can be used to make explicit predictions about how visual perception is modulated by different stimuli and tasks and thus could help contribute this debate.

In our model, the agent's attentional state is altered slowly, based on feedback on its actions on many trials. However, in reality, attention can

be quickly redirected following explicit sensory cues or instructions. Our model could be extended to account for these quick changes in attentional state by including additional variables in the internal model to represent the current behavioral context (e.g., the location of the detection target). Thus, on any given trial, the agent would first have to infer the behavioral context based on all available sources of information (e.g., received rewards, sensory cues, instructions, or prior experience in the task). The inferred context would then determine the attentional state that optimized the internal sensory representation toward the task. Such a model could be used to predict how people's attentional state is altered in real time as a result of newly received information. Our goal, however, was more modest: we sought to explore the effects (rather than the temporal dynamics; Yu et al., 2009) of optimizing the internal sensory representation towards a given task.

In this letter, we put forward a very general framework for predicting how task demands alter visual processing. We then showed that given certain assumptions about the internal model and behavioral task, this framework predicts attention-dependent changes to neural responses that are consistent with existing phenomenological models of attention. However, although the assumptions of our model are based on functional principles, in order to truly derive the effects of attention, it would be desirable to construct a more sophisticated model of natural images in which model parameters are learned directly from natural image statistics (as opposed to artificial data). In the past, this approach has been highly successful in understanding the passive properties of visual neurons. In the future, it could be used to make quantitative and testable predictions about how different behavioral tasks alter visual processing and perception.

## Acknowledgments

---

We thank Odelia Schwartz, Ruben Coen-Caglie, and Peter Dayan for their helpful comments and feedback on an earlier version of this letter. This research was supported by funding from the Engineering and Physical Sciences Research Council and the Medical Research Council of Great Britain.

## References

---

- Berkes, P., Turner, R., & Sahani, M. (2008). On sparsity and overcompleteness in image models. In J. C. Platt, D. Koller, Y. Singer, & S. Roweis (Eds.), *Advances in neural information processing systems*, 22. Cambridge, MA: MIT Press.
- Chalk, M., Seitz, A. R., & Seriès, P. (2010). Rapidly learned stimulus expectations alter perception of motion. *Journal of Vision*, 10(8), 2.
- Chikkerur, S., Serre, T., Tan, C., & Poggio, T. (2010). What and where: A Bayesian inference theory of attention. *Vision Research*, 50(22), 2233–2247.
- Dayan, P., & Daw, N. D. (2008). Decision theory, reinforcement learning, and the brain. *Cognitive, Affective, and Behavioral Neuroscience*, 8(4), 429–453.

- Dayan, P., & Solomon, J. A. (2010). Selective Bayes: Attentional load and crowding. *Vision Research*, *50*(22), 2248–2260.
- Dayan, P., & Zemel, R. S. (1999). Statistical models and sensory attention. In *Proceedings of the International Conference on Artificial Neural Networks* (Vol. 2, pp. 1017–1022). Piscataway, NJ: IEEE.
- Deneve, S. (2008). Bayesian spiking neurons I: Inference. *Neural Computation*, *20*(1), 91–117.
- Downing, C. J. (1988). Expectancy and visual-spatial attention: Effects on perceptual quality. *Journal of Experimental Psychology Human Perception and Performance*, *14*(2), 188–202.
- Eckstein, M. P., Abbey, C. K., Pham, B. T., & Shimozaki, S. S. (2004). Perceptual learning through optimization of attentional weighting: Human versus optimal Bayesian learner. *Journal of Vision*, *4*(12), 1006–1019.
- Fahle, M. (2005). Perceptual learning: Specificity versus generalization. *Current Opinion in Neurobiology*, *15*(2), 154–160.
- Fawcett, T. (2006). An introduction to ROC analysis. *Pattern Recognition Letters*, *27*, 861–874.
- Fiser, J., Berkes, P., Orbán, G., & Lengyel, M. (2010). Statistically optimal perception and learning: From behavior to neural representations. *Trends in Cognitive Sciences*, *14*(3), 119–130.
- Ghose, G. M. (2009). Attentional modulation of visual responses by flexible input gain. *Journal of Neurophysiology*, *101*(4), 2089–2106.
- Hyvärinen, A. (2010). Statistical models of natural images and cortical visual representation. *Topics in Cognitive Science*, *2*(2), 251–264.
- Jazayeri, M. (2007). Integration of sensory evidence in motion discrimination. *Journal of Vision*, *7*(12), 1–7.
- Jiang, Y., & Chun, M. M. (2001). Selective attention modulates implicit learning. *Quarterly Journal of Experimental Psychology*, *54*(4), 1105–1124.
- Karklin, Y., & Lewicki, M. S. (2003). Learning higher-order structures in natural images. *Network (Bristol, England)*, *14*(3), 483–499.
- Lee, J., & Maunsell, J. H. R. (2009). A normalization model of attentional modulation of single unit responses. *PloS One*, *4*(2), e4651.
- Liu, Y., Yu, A., & Holmes, P. (2009). Dynamical analysis of Bayesian inference models for the Eriksen task. *Neural Computation*, *21*, 1520–1553.
- Luck, S. J., Chelazzi, L., Hillyard, S. A., & Desimone, R. (1997). Neural mechanisms of spatial selective attention in areas V1, V2, and V4 of macaque visual cortex. *Journal of Neurophysiology*, *77*, 24–42.
- Ma, W. J., Beck, J. M., Latham, P. E., & Pouget, A. (2006). Bayesian inference with probabilistic population codes. *Nature Neuroscience*, *9*(11), 1432–1438.
- Martinez-Trujillo, J., & Treue, S. (2004). Feature-based attention increases the selectivity of population responses in primate visual cortex. *Current Biology*, *14*, 744–751.
- McAdams, C. J., & Maunsell, J. H. R. (1999). Effects of attention on orientation-tuning functions of single neurons in macaque cortical area V4. *Journal of Neuroscience*, *19*(1), 431–441.
- Moran, J., & Desimone, R. (1985). Selective attention gates visual processing in the extrastriate cortex. *Science*, *229*(4715), 782–784.

- Olshausen, B. A., & Field, D. J. (1996). Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, *381*(6583), 607–609.
- Olshausen, B. A., & Field, D. J. (1997). Sparse coding with an overcomplete basis set: A strategy employed by V1? *Vision Research*, *37*(23), 3311–3325.
- Pestilli, F., & Carrasco, M. (2005). Attention enhances contrast sensitivity at cued and impairs it at uncued locations. *Vision Research*, *45*(14), 1867–1875.
- Posner, M. I., Snyder, C. R. R., & Davidson, B. J. (1980). Attention and the detection of signals. *Journal of Experimental Psychology*, *109*(2), 160–174.
- Puertas, G., Bornschein, J., & Lücke, J. (2010). The maximal causes of natural scenes are edge filters. In J. Lafferty, C.K.I. Williams, J. Shawe-Taylor, R. S. Zemel, & A. Culotta (Eds.), *Advances in neural information processing*, 23 (pp. 1939–1947). Red Hook, NY: Curran.
- Rao, R.P.N. (2005). Bayesian inference and attentional modulation in the visual cortex. *NeuroReport*, *16*(16), 1843–1848.
- Reynolds, J. H., & Chelazzi, L. (2004). Attentional modulation of visual processing. *Annual Review of Neuroscience*, *27*, 611–647.
- Reynolds, J. H., Chelazzi, L., & Desimone, R. (1999). Competitive mechanisms subserve attention in macaque areas V2 and V4. *Journal of Neuroscience*, *19*(5), 1736–1753.
- Reynolds, J. H., & Heeger, D. J. (2009). The normalization model of attention. *Neuron*, *61*(2), 168–185.
- Reynolds, J. H., Pasternak, T., & Desimone, R. (2000). Attention increases sensitivity of V4 neurons. *Neuron*, *26*(3), 703–714.
- Sahani, M. (2004). A biologically plausible algorithm for reinforcement-shaped representational learning. In S. Thrun, L. Saul, & B. Scholköpfung (Eds.), *Advances in neural information processing*, 16 (pp. 1287–1294). Cambridge, MA: MIT Press.
- Schwartz, O., & Coen-Cagli, R. (2013). Visual attention and flexible normalization pools. *Journal of Vision*, *13*(1), 1–24.
- Seitz, A. R., Kim, D., & Watanabe, T. (2009). Rewards evoke learning of unconsciously processed visual stimuli in adult humans. *Neuron*, *61*(5), 700–707.
- Seriès, P., Lorenceau, J., & Frégnac, Y. (2003). The silent surround of V1 receptive fields: Theory and experiments. *Journal of Physiology*, *97*(4–6), 453–474.
- Shelton, J. A., Bornschein, J., Sheikh, A. S., Berkes, P., & Lücke, J. (2011). Select and sample: A model of efficient neural inference and learning. In J. Shawe-Taylor, R. S. Zemel, P. Bartlett, F. Pereira, & K. Q. Weinberger (Eds.), *Advances in neural information processing systems*, 24. Red Hook, NY: Curran.
- Sundberg, K. A., Mitchell, J. F., & Reynolds, J. H. (2009). Spatial attention modulates center-surround interactions in macaque visual area V4. *Neuron*, *61*(6), 952–963.
- Whiteley, L., & Sahani, M. (2012, June). Attention in a bayesian framework. *Frontiers in Human Neuroscience*, *6*, 100.
- Williford, T., & Maunsell, J. H. R. (2006). Effects of spatial attention on contrast response functions in macaque area V4. *Journal of Neurophysiology*, *96*(1), 40–54.
- Womelsdorf, T., Anton-Erxleben, K., Pieper, F., & Treue, S. (2006). Dynamic shifts of visual receptive fields in cortical area MT by spatial attention. *Nature Neuroscience*, *9*(9), 1156–1160.

Yu, A., & Dayan, P. (2005). Uncertainty, neuromodulation, & attention. *Neuron*, *46*(4), 681–692.

Yu, A. J., Dayan, P., & Cohen, J. D. (2009). Dynamics of attentional selection under conflict: Toward a rational Bayesian account. *Journal of Experimental Psychology: Human Perception and Performance*, *35*(3), 700–717.

---

Received November 1, 2012; accepted April 9, 2013.