

A. Extraction of early visual features

With r , g and b being the red, green and blue channels of the input image, an intensity image I is obtained as $I = (r+g+b)/3$. I is used to create a Gaussian pyramid $I(\sigma)$, where $\sigma \in [0..8]$ is the scale. The r , g and b channels are normalized by I in order to decouple hue from intensity. However, because hue variations are not perceivable at very low luminance (and hence are not salient), normalization is only applied at the locations where I is larger than $1/10$ of its maximum over the entire image (other locations yield zero r, g and b). Four broadly-tuned color channels are created: $R = r - (g + b)/2$ for red, $G = g - (r + b)/2$ for green, $B = b - (r + g)/2$ for blue, and $Y = (r + g)/2 - |r - g|/2 - b$ for yellow (negative values are set to zero). Four Gaussian pyramids $R(\sigma), G(\sigma), B(\sigma)$ and $Y(\sigma)$ are created from these color channels.

Center-surround differences (\ominus defined previously) between a “center” fine scale c and a “surround” coarser scale s yield the feature maps. The first set of feature maps is concerned with intensity contrast, which in mammals is detected by neurons sensitive either to dark centers on bright surrounds, or to bright centers on dark surrounds [12]. Here, both types of sensitivities are simultaneously computed (using a rectification) in a set of six maps $\mathcal{I}(c, s)$, with $c \in \{2, 3, 4\}$ and $s = c + \delta$, $\delta \in \{3, 4\}$:

$$\mathcal{I}(c, s) = |I(c) \ominus I(s)| \quad (1)$$

A second set of maps is similarly constructed for the color channels, which in cortex are represented using a so-called “color double-opponent” system: In the center of their receptive field, neurons are excited by one color (e.g., red) and inhibited by another (e.g., green), while the converse is true in the surround. Such spatial and chromatic opponency exists for the red/green, green/red, blue/yellow and yellow/blue color pairs in human primary visual cortex [13]. Accordingly, maps $\mathcal{RG}(c, s)$ are created in the model to simultaneously account for red/green and green/red double opponency (**Eq. 2**), and $\mathcal{BY}(c, s)$ for blue/yellow and yellow/blue double opponency (**Eq. 3**):

$$\mathcal{RG}(c, s) = |(R(c) - G(c)) \ominus (G(s) - R(s))| \quad (2)$$

$$\mathcal{BY}(c, s) = |(B(c) - Y(c)) \ominus (Y(s) - B(s))| \quad (3)$$

Local orientation information is obtained from I using oriented Gabor pyramids $O(\sigma, \theta)$, where $\sigma \in [0..8]$ represents the scale and $\theta \in \{0^\circ, 45^\circ, 90^\circ, 135^\circ\}$ is the preferred orientation [11]. (Gabor filters, which are the product of a cosine grating and a 2D Gaussian envelope, approximate the receptive field sensitivity profile (impulse response) of orientation-selective neurons in primary visual cortex [12].) Orientation feature maps, $\mathcal{O}(c, s, \theta)$, encode, as a group, local orientation contrast between the center and surround scales:

$$\mathcal{O}(c, s, \theta) = |O(c, \theta) \ominus O(s, \theta)| \quad (4)$$

In total, 42 feature maps are computed: Six for intensity, 12 for color and 24 for orientation.

B. The Saliency Map

The purpose of the saliency map is to represent the conspicuity — or “saliency” — at every location in the visual field by a scalar quantity, and to guide the selection of attended locations, based on the spatial distribution of saliency. A combination of the feature maps provides bottom-up input to the saliency map, modeled as a dynamical neural network.

One difficulty in combining different feature maps is that they represent *a priori* not comparable modalities, with different dynamic ranges and extraction mechanisms. Also, because all 42

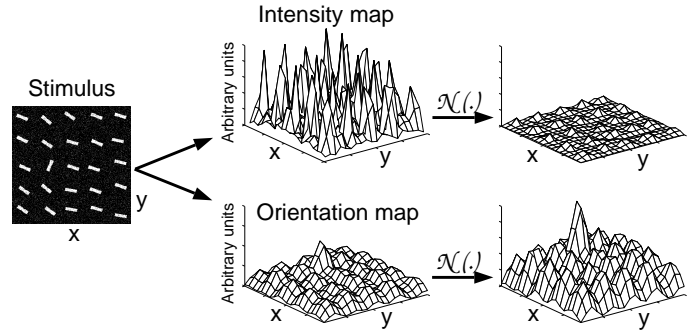


Fig. 2. The normalization operator $\mathcal{N}(\cdot)$.

feature maps are combined, salient objects appearing strongly in only a few maps may be masked by noise or less salient objects present in a larger number of maps.

In the absence of top-down supervision, we propose a map normalization operator, $\mathcal{N}(\cdot)$, which globally promotes maps in which a small number of strong peaks of activity (conspicuous locations) is present, while globally suppressing maps which contain numerous comparable responses. $\mathcal{N}(\cdot)$ consists of (**Fig. 2**): 1) Normalizing the values in the map to a fixed range $[0..M]$, in order to eliminate modality-dependent amplitude differences; 2) finding the location of the map’s global maximum M and computing the average \bar{m} of all its other local maxima; 3) globally multiplying the map by $(M - \bar{m})^2$.

Only local maxima of activity are considered such that $\mathcal{N}(\cdot)$ compares responses associated with meaningful “activation spots” in the map and ignores homogenous areas. Comparing the maximum activity in the entire map to the average over all activation spots measures how different the most active location is from the average. When this difference is large, the most active location stands out, and we strongly promote the map. When the difference is small, the map contains nothing unique and is suppressed. The biological motivation behind the design of $\mathcal{N}(\cdot)$ is that it coarsely replicates cortical lateral inhibition mechanisms, in which neighboring similar features inhibit each other *via* specific, anatomically-defined connections [15].

Feature maps are combined into three “conspicuity maps”, $\bar{\mathcal{I}}$ for intensity (**Eq. 5**), $\bar{\mathcal{C}}$ for color (**Eq. 6**), and $\bar{\mathcal{O}}$ orientation (**Eq. 7**), at the scale ($\sigma = 4$) of the saliency map. They are obtained through across-scale addition, “ \oplus ”, which consists of reduction of each map to scale 4 and point-by-point addition:

$$\bar{\mathcal{I}} = \bigoplus_{c=2}^4 \bigoplus_{s=c+3}^{c+4} \mathcal{N}(\mathcal{I}(c, s)) \quad (5)$$

$$\bar{\mathcal{C}} = \bigoplus_{c=2}^4 \bigoplus_{s=c+3}^{c+4} [\mathcal{N}(\mathcal{RG}(c, s)) + \mathcal{N}(\mathcal{BY}(c, s))] \quad (6)$$

For orientation, four intermediary maps are first created by combination of the six feature maps for a given θ , and are then combined into a single orientation conspicuity map:

$$\bar{\mathcal{O}} = \sum_{\theta \in \{0^\circ, 45^\circ, 90^\circ, 135^\circ\}} \mathcal{N} \left(\bigoplus_{c=2}^4 \bigoplus_{s=c+3}^{c+4} \mathcal{N}(\mathcal{O}(c, s, \theta)) \right) \quad (7)$$

The motivation for the creation of three separate channels, $\bar{\mathcal{I}}$, $\bar{\mathcal{C}}$ and $\bar{\mathcal{O}}$, and their individual normalization is the hypothesis that similar features compete strongly for saliency, while different modalities contribute independently to the saliency map.

The three conspicuity maps are normalized and summed into the final input \mathcal{S} to the saliency map:

$$\mathcal{S} = \frac{1}{3} (\mathcal{N}(\bar{\mathcal{I}}) + \mathcal{N}(\bar{\mathcal{C}}) + \mathcal{N}(\bar{\mathcal{O}})) \quad (8)$$

At any given time, the maximum of the saliency map (SM) defines the most salient image location, to which the focus of attention (FOA) should be directed. We could now simply select the most active location as defining the point where the model should next attend to. However, in a neuronally-plausible implementation, we model the SM as a 2D layer of leaky *integrate-and-fire* neurons at scale 4. These model neurons consist of a single capacitance which integrates the charge delivered by synaptic input, of a leakage conductance, and of a voltage threshold. When threshold is reached, a prototypical spike is generated, and the capacitive charge is shunted to zero [14]. The SM feeds into a biologically-plausible 2D “winner-take-all” (WTA) neural network [4], [1] at scale $\sigma = 4$, in which synaptic interactions among units ensure that only the most active location remains, while all other locations are suppressed.

The neurons in the SM receive excitatory inputs from \mathcal{S} and are all independent. The potential of SM neurons at more salient locations hence increases faster (these neurons are used as pure integrators and do not fire). Each SM neuron excites its corresponding WTA neuron. All WTA neurons also evolve independently of each other, until one (the “winner”) first reaches threshold and fires. This triggers three simultaneous mechanisms (Fig. 3): 1) The FOA is shifted to the location of the winner neuron; 2) the global inhibition of the WTA is triggered and completely inhibits (resets) all WTA neurons; 3) local inhibition is transiently activated in the SM, in an area with the size and new location of the FOA; this not only yields dynamical shifts of the FOA, by allowing the next most salient location to subsequently become the winner, but it also prevents the FOA from immediately returning to a previously attended location. Such an “inhibition of return” has been demonstrated in human visual psychophysics [16]. In order to slightly bias the model to subsequently jump to salient locations spatially close to the currently attended location, a small excitation is transiently activated in the SM, in a near surround of the FOA (“proximity preference” rule of Koch and Ullman [4]).

Since we do not model any top-down attentional component, the FOA is a simple disk whose radius is fixed to one sixth of the smaller of the input image width or height. The time constants, conductances, and firing thresholds of the simulated neurons were chosen (see ref. [17] for details) so that the FOA jumps from one salient location to the next in approximately 30–70 ms (simulated time), and that an attended area is inhibited for approximately 500–900 ms (Fig. 3), as has been observed psychophysically [16]. The difference in the relative magnitude of these delays proved sufficient to ensure thorough scanning of the image, and prevented cycling through only a limited number of locations. All parameters are fixed in our implementation [17], and the system proved stable in time for all images studied.

C. Comparison with spatial frequency content models

Reinagel and Zador [18] recently used an eye-tracking device to analyze the local spatial frequency distributions along eye scan paths generated by humans while free-viewing grayscale images. They found the spatial frequency content at the fixated locations to be significantly higher than, on average, at random locations. Although eye trajectories can differ from attentional trajectories under volitional control [1], visual attention is often thought as a pre-oculomotor mechanism, strongly influencing

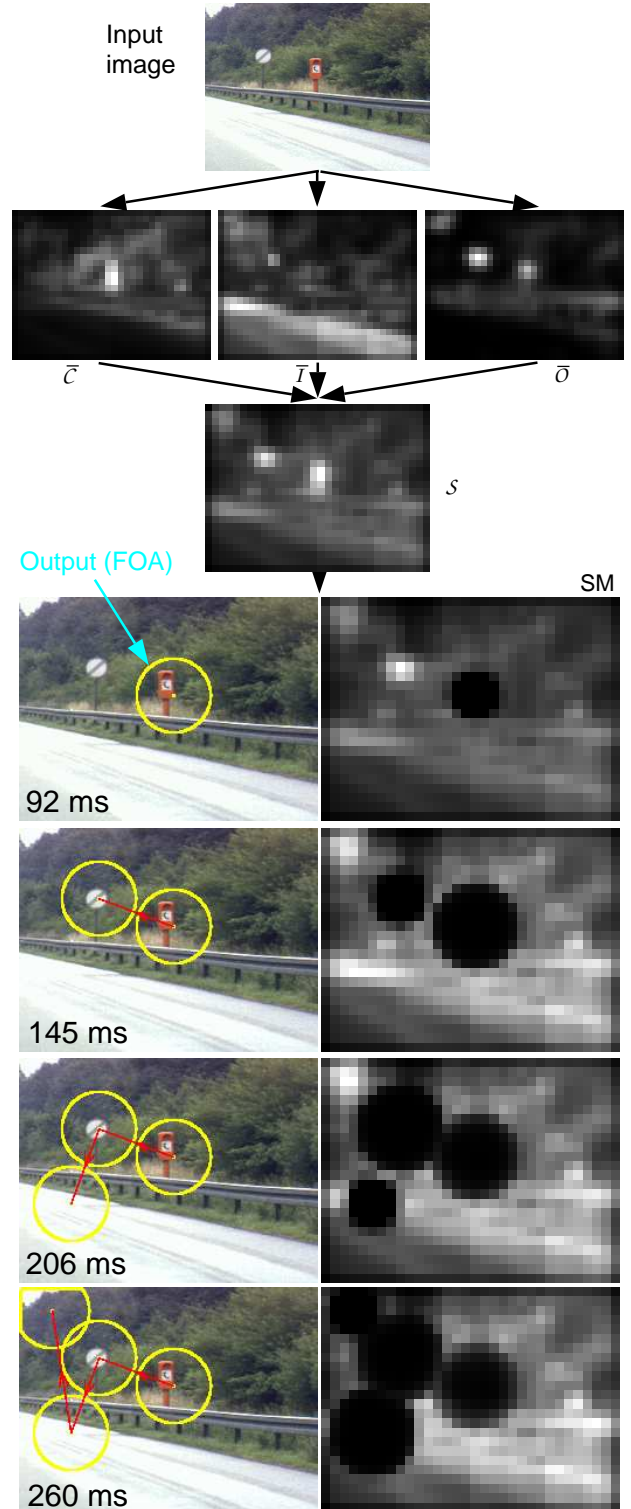


Fig. 3. Example of operation of the model with a natural image. Parallel feature extraction yields the three conspicuity maps for color contrasts ($\bar{\mathcal{C}}$), intensity contrasts ($\bar{\mathcal{I}}$), and orientation contrasts ($\bar{\mathcal{O}}$). These are combined to form input \mathcal{S} to the saliency map (SM). The most salient location is the orange telephone box, which appeared very strongly in $\bar{\mathcal{C}}$; it becomes the first attended location (92 ms simulated time). After the inhibition-of-return feedback inhibits this location in the saliency map, the next most salient locations are successively selected.

free-viewing. It was hence interesting to investigate whether our model would reproduce the findings of Reinagel and Zador.

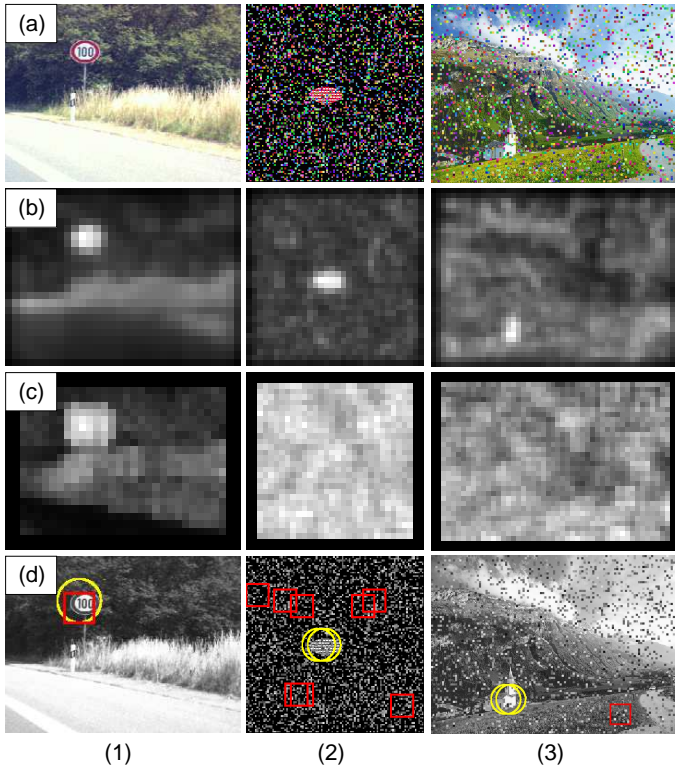


Fig. 4. Examples of color images (a), the corresponding saliency map inputs (b), spatial frequency content (SFC) maps (c), locations at which input to the saliency map was higher than 98% of its maximum (d; yellow circles), and image patches for which the SFC was higher than 98% of its maximum (d; red squares). The saliency maps are very robust to noise, while SFC is not.

We constructed a simple measure of spatial frequency content (SFC): At a given image location, a 16×16 image patch is extracted from each $I(2)$, $R(2)$, $G(2)$, $B(2)$ and $Y(2)$ map, and 2D Fast Fourier Transforms (FFTs) are applied to the patches. For each patch, a threshold is applied to compute the number of non-negligible FFT coefficients; the threshold corresponds to the FFT amplitude of a just perceivable grating (1% contrast). The SFC measure is the average of the numbers of non-negligible coefficients in the five corresponding patches. The size and scale of the patches were chosen such that the SFC measure is sensitive to approximately the same frequency and resolution ranges as our model; also, our SFC measure is computed in the RGB channels as well as in intensity, like the model. Using this measure, an SFC map is created at scale 4 for comparison with the saliency map (Fig. 4).

III. RESULTS AND DISCUSSION

Although the concept of a saliency map has been widely used in focus-of-attention models [1], [3], [7], little detail is usually provided about its construction and dynamics. Here we examine how the feedforward feature extraction stages, the map combination strategy, and the temporal properties of the saliency map all contribute to the overall system performance.

A. General performance

The model was extensively tested with artificial images to ensure proper functioning. For example, several objects of same shape but varying contrast with the background were attended to in order of decreasing contrast. The model proved very robust to the addition of noise to such images (Fig. 5), particularly

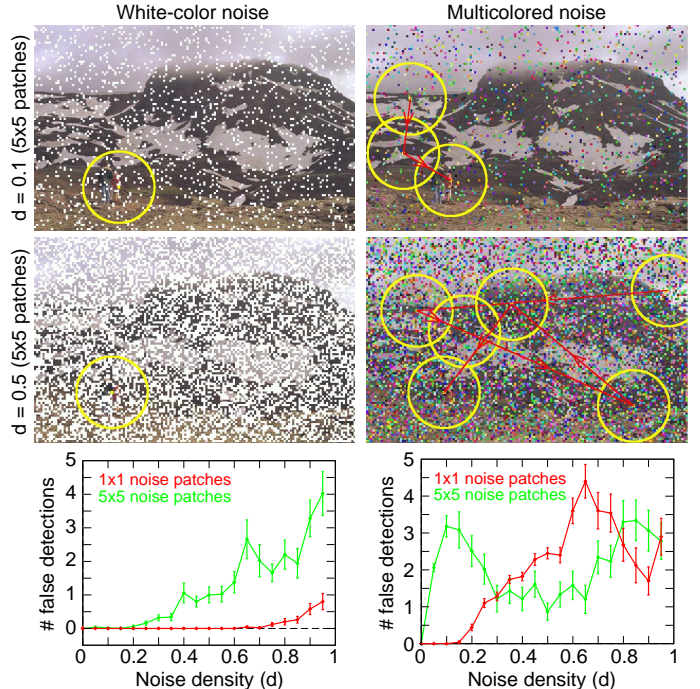


Fig. 5. Influence of noise on detection performance, illustrated with a 768×512 scene in which a target (two people) is salient by its unique color contrast. The mean \pm S.E. of false detections before target found is shown as a function of noise density for 50 instantiations of the noise. The system is very robust to noise which does not directly interfere with the main feature of the target (left; intensity noise and color target). When the noise has similar properties to the target, it impairs the target's saliency and the system first attends to objects salient for other features (here, coarse-scale variations of intensity).

if the properties of the noise (e.g., its color) were not directly conflicting with the main feature of the target.

The model was able to reproduce human performance for a number of pop-out tasks [7], using images of the type shown in Fig. 2. When a target differed from an array of surrounding distractors by its unique orientation (like in Fig. 2), color, intensity or size, it was always the first attended location, irrespectively of the number of distractors. Contrarily, when the target differed from the distractors only by a conjunction of features (e.g., it was the only red-horizontal bar in a mixed array of red-vertical and green-horizontal bars), the search time necessary to find the target increased linearly with the number of distractors. Both results have been widely observed in humans [7], and are discussed in Section III-B.

We also tested the model with real images, ranging from natural outdoor scenes to artistic paintings, and using $\mathcal{N}(\cdot)$ to normalize the feature maps (Fig. 3 and ref. [17]). With many such images, it is difficult to objectively evaluate the model, because no objective reference is available for comparison, and observers may disagree on which locations are the most salient. However, in all images studied, most attended locations were objects of interest, such as faces, flags, persons, buildings or vehicles.

Model predictions were compared to the measure of local SFC, in an experiment similar to that of Reinagel and Zador [18], using natural scenes with salient traffic signs (90 images), red soda can (104 images), or vehicle's emergency triangle (64 images). Similar to Reinagel and Zador's findings, the SFC at attended locations was significantly higher than the average SFC, by a factor decreasing from 2.5 ± 0.05 at the first attended location to 1.6 ± 0.05 at the 8th attended location. Although

this result does not necessarily indicate similarity between human eye fixations and the model's attentional trajectories, it indicates that the model, like humans, is attracted to "informative" image locations, according to the common assumption that regions with richer spectral content are more informative. The SFC map was similar to the saliency map for most images (e.g., Fig. 4.1). However, both maps differed substantially for images with strong, extended variations of illumination or color (e.g., due to speckle noise): While such areas exhibited uniformly high SFC, they had low saliency because of their uniformity (Figs. 4.2, 4.3). In such images, the saliency map was usually in better agreement with our subjective perception of saliency. Quantitatively, for the 258 images studied here, the SFC at attended locations was significantly lower than the maximum SFC, by a factor decreasing from 0.90 ± 0.02 at the first attended location to 0.55 ± 0.05 at the 8th attended location: While the model was attending to locations with high SFC, these were not necessarily the locations with highest SFC. It consequently seems that saliency is more than just a measure of local SFC. The model, which implements within-feature spatial competition captured subjective saliency better than the purely local SFC measure.

B. Strengths and limitations

We have proposed a model whose architecture and components mimic the properties of primate early vision. Despite its simple architecture and feedforward feature extraction mechanisms, the model is capable of strong performance with complex natural scenes. For example, it quickly detected salient traffic signs of varied shapes (round, triangular, square, rectangular), colors (red, blue, white, orange, black) and textures (letter markings, arrows, stripes, circles), although it had not been designed for this purpose. Such strong performance reinforces the idea that a unique saliency map, receiving input from early visual processes, could effectively guide bottom-up attention in primates [4], [10], [5], [8]. From a computational viewpoint, the major strength of this approach lies in the massively parallel implementation, not only of the computationally expensive early feature extraction stages, but also of the attention focusing system. More than previous models based extensively on relaxation techniques [5], our architecture could easily allow for real-time operation on dedicated hardware.

The type of performance which can be expected from this model critically depends on one factor: Only object features explicitly represented in at least one of the feature maps can lead to pop-out, that is, rapid detection independently of the number of distracting objects [7]. Without modifying the pre-attentive feature extraction stages, our model cannot detect conjunctions of features. While our system immediately detects a target which differs from surrounding distractors by its unique size, intensity, color or orientation (properties which we have implemented because they have been very well characterized in primary visual cortex), it will fail at detecting targets salient for unimplemented feature types (e.g., T junctions or line terminators, for which the existence of specific neural detectors remains controversial). For simplicity, we also have not implemented any recurrent mechanism within the feature maps, and hence cannot reproduce phenomena like contour completion and closure, important for certain types of human pop-out [19]. In addition, at present our model does not include any magnocellular motion channel, known to play a strong role in human saliency [5].

A critical model component is the normalization $\mathcal{N}(\cdot)$, which provided a general mechanism for computing saliency in any situation. The resulting saliency measure implemented by the

model, although often related to local SFC, was closer to human saliency because it implemented spatial competition between salient locations. Our feed-forward implementation of $\mathcal{N}(\cdot)$ is faster and simpler than previously proposed iterative schemes [5]. Neuronally, spatial competition effects similar to $\mathcal{N}(\cdot)$ have been observed in the non-classical receptive field of cells in striate and extrastriate cortex [15].

In conclusion, we have presented a conceptually simple computational model for saliency-driven focal visual attention. The biological insight guiding its architecture proved efficient in reproducing some of the performances of primate visual systems. The efficiency of this approach for target detection critically depends on the features types implemented. The framework presented here can consequently be easily tailored to arbitrary tasks through the implementation of dedicated feature maps.

ACKNOWLEDGMENTS

We thank Werner Ritter and Daimler-Benz for the traffic sign images, Pietro Perona and both reviewers for excellent suggestions. Supported by the National Science Foundation and the Office of Naval Research.

REFERENCES

- [1] J. K. Tsotsos, S. M. Culhane, W. Y. K. Wai, Y. H. Lai, N. Davis, F. Nuflo, "Modelling visual attention via selective tuning," *Artificial Intelligence*, vol. 78, no. 1-2, pp. 507-545, Oct. 1995.
- [2] E. Niebur and C. Koch, "Computational architectures for attention," R. Parasuraman, (Ed.), *The attentive brain*, Cambridge, MA:MIT Press, pp. 163-186, 1998.
- [3] B.A. Olshausen, C.H. Anderson CH and D.C. Van Essen, "A neurobiological model of visual attention and invariant pattern recognition based on dynamic routing of information," *J Neuroscience*, vol. 13, no. 11, pp. 4700-4719, Nov. 1993.
- [4] C. Koch and S. Ullman, "Shifts in selective visual attention: towards the underlying neural circuitry," *Human Neurobiology*, vol. 4, pp. 219-227, 1985.
- [5] R. Milanese, S. Gil and T. Pun, "Attentive Mechanisms for Dynamic and Static Scene Analysis," *Optical Engineering*, vol. 34, no. 8, pp.2428-2434, Aug. 1995.
- [6] S. Baluja and D.A. Pomerleau, "Expectation-based Selective Attention for Visual Monitoring and Control of a Robot Vehicle," *Robotics and Autonomous Systems*, vol. 22, no. 3-4, pp. 329-344, Dec. 1997.
- [7] A.M. Treisman and G. Gelade, "A feature-integration theory of attention," *Cognitive Psychology*, vol. 12, no. 1, pp. 97-136, Jan. 1980.
- [8] J.P. Gottlieb, M. Kusunoki and M.E. Goldberg, "The representation of visual salience in monkey parietal cortex," *Nature*, vol. 391, no. 6666, pp. 481-484, Jan. 1998.
- [9] D.L. Robinson, S.E. Peterson, "The pulvinar and visual salience," *Trends in Neurosciences*, vol. 15, no. 4, pp. 127-132, Apr. 1992.
- [10] J.M. Wolfe, "Guided search 2.0: a revised model of visual search," *Psychonomic Bulletin Review*, vol. 1, pp. 202-238, 1994.
- [11] H. Greenspan, S. Belongie, R. Goodman, P. Perona, S. Rakshit, C. H. Anderson, "Overcomplete steerable pyramid filters and rotation invariance," *Proc. IEEE Computer Vision and Pattern Recognition (CVPR)*, Seattle, Washington, pp. 222-228, Jun. 1994.
- [12] A.G. Leventhal, *The Neural Basis of Visual Function (Vision and Visual Dysfunction Vol. 4)*, Boca Raton, FL:CRC Press, 1991.
- [13] S. Engel, X. Zhang and B. Wandell, "Colour tuning in human visual cortex measured with functional magnetic resonance imaging," *Nature*, vol. 388, no. 6637, pp. 68-71, Jul. 1997.
- [14] C. Koch, "Biophysics of Computation: Information Processing in Single Neurons," Oxford University Press (in press).
- [15] M. W. Cannon and S. C. Fullenkamp, "A Model for Inhibitory Lateral Interaction Effects in Perceived Contrast," *Vision Research*, vol. 36, no. 8, pp. 1115-1125, Apr. 1996.
- [16] M. I. Posner and Y. Cohen, "Components of visual orienting", *Attention and Performance X*, (H. Bouma, D.G. Bouwhuis eds), Hillsdale, NJ:L. Erlbaum, pp. 531-556, 1984.
- [17] The C++ implementation of the model and numerous examples of attentional predictions on natural and synthetic images can be retrieved from <http://www.klab.caltech.edu/~itti/attention/>
- [18] P. Reinagel and A.M. Zador, "The Effect of Gaze on Natural Scene Statistics," *Neural information and coding workshop*, Snowbird, Utah, 16-20 Mar. 1997.
- [19] I. Kovacs and B. Julesz, "A closed curve is much more than an incomplete one: effect of closure in figure-ground segmentation," *Proc. Nat'l Academy of Sciences, U.S.A.*, vol. 90, no. 16, pp. 7495-7497, Aug. 1993.