

Saliency map based attention control for the RoboCup SPL

Juan F. García, Francisco J. Rodríguez, Camino Fernández, Vicente Matellán

Grupo de Robótica - Escuela de Ingenierías Industrial e Informática

Universidad de León

{jfgarsl, frodl, camino.fernandez, vicente.matellan}@unileon.es

Abstract

Attention mechanisms can be used both for reducing the amount of perceptual information to be processed and for restricting all available actions to only those useful for a given scenario. Information reduction improves performance and action restriction allows for a more precise interaction with our environment. In this paper we present the design of an attention control mechanism based on a saliency map and its implementation in the SPL's Nao robot. The results obtained are analysed and future works derived from that analysis are presented.

Index terms – RoboCup, attention, saliency, map, humanoid

1 Introduction

Attention is a natural tool which allows animals to locate relevant objects or areas in a given scene, discarding the rest of elements present and thus reducing the amount of information to deal with [2, 4]. The areas marked as conspicuous and the objects they contain restrict all our possible actions to those which can be specifically applied to them, discarding other distracting elements [15, 17].

Vision and control systems in Robotics are usually implemented in an impulse-analysis-response fashion. Given a visual impulse, the analysis subsystem generates a “world model” which is then used by the response module to generate an action. In this case, vision is just a step previous to planning. However, attention can be used to further relate these two sys-

tems: control system can establish the kind of objects that should be looked for (top-down, control modulates attention) and attended locations restrict what can be done in that moment (bottom-up, attention modulates control) [3].

The latest attention models are mostly bioinspired and try to reproduce the way primates' and humans' attention works [6, 8, 19]. Color contrast, intensity difference, orientation and motion are just some of the key elements considered by these models.

In this paper we present a bioinspired attention model mainly based on Itti et al. research [8, 9, 10] which falls under the bottom-up attention category. It has been developed according to the Standard Platform League (SPL) regulations and is intended to be tested during future Robocup¹ events in 2011.

For a better performance and adaptation to this environment, our model does not use all the maps the original utilises. Input image size is also reduced by means of a virtual fovea mask, further releasing computational resources. The use of this attention system will also allow us to participate in some of the latest proposed challenges, like, for instance, the “any ball” challenge, with better results than classic filter and segmentation algorithms provide.

The rest of the paper is organised as follows. In the second section, some of the most notable attention models are enumerated. In the third section, the attention model is explained, both the principles and the software structure are detailed. In the fourth section,

¹<http://www.robocup.org/>

the attention algorithms used are described. In the fifth section, experiments used for the model validation are summarised. Finally, in the last section, the results obtained are discussed and also the future works envisioned are enumerated.

2 Attention models

Animals, and humans specifically, can change their focus of attention either by moving their fixation point across the visual scene or by focusing on a given area of the current visual field. The former is known as “overt attention” and the latter, which is the one we mainly describe in this article, as “covert attention” [23]. Covert changes are much faster (up to five times) than overt ones, which makes this early attention an important tool to decide whether it is suitable or not to change the current fixation point (move our eyes or even the head).

Several attention models have been proposed over the years, mainly from a psychological and neurological point of view [12, 14]. Natural attention is the starting point of all of them. Since a detailed analysis transcends the scope of this paper, a list of those more related to this work is given:

- Classic Attention model by Koch and Ullman [13]. Several feature maps are extracted from the input image and then used to build a saliency map. A WTS (*winner-takes-all*) process will then select the more relevant areas in this map and direct attention to them. It is the base of most of the other models explained in here.
- Wolfe’s Guided Search model [22]. Based on Koch and Ullman model, it starts with the computation of basic features, such as color and orientation, which are then used to build the so called feature maps. These maps are finally merged in an activation map which will be used for guiding the attention to the most relevant areas (those with higher values in the map).
- Saliency Map models. Itti et al. [8, 9, 10, 11] developed a model closely related

to Koch and Ullman studies. This model builds up a saliency map to guide attention using color, intensity, orientation and movement maps which are extracted from the input images.

All these models are often called “Feature-Based Attention Models”. Their main objective is identifying the more conspicuous areas in the current scene. There are several other approaches created to model attention, such as “Connectionist Attention Models” [5, 7, 16], which are oriented to create a reference frame for specific objects or some of their environmental interaction features (for instance, specific movement patterns).

The model described in this paper is an adaptation to the Robocup SPL environment of Itti’s proposal. While being conceptually simple, it offers great results while not consuming a high amount of resources. To further prioritise performance, several of its elements have been simplified: some of its maps are dispensed while the use of a fovea mask reduces original images size. The model will be further reviewed in the next section.

3 Model

Our model is based on Itti et al. saliency map attention model [8, 9, 10, 11]. At any given time, the maximum registered in the saliency defines the most important region from an attentive point of view.

To build up the saliency map in our model, two maps are used: an intensity map and a color map. The other two maps of the original model (orientation and movement) are dispensed since we do not find them necessary for our environment.

The maps assign high values to those areas which stand out in the magnitude they measure: intensity map will assign high values to those areas the intensity (light) of which changes a lot in relation to their surround, while color maps will do the same for the ones with a high color contrast.

The maps are obtained using the original camera image. To reduce the amount of pixels

to be computed, and thus improving performance, a virtual fovea mask which simulates the human eye progressive resolution decrement is used. The further a region is from the center of the image, the greater the amount of masked pixels will be (masked pixels will not be analysed). Section 4.1 shows a more detailed description of this fovea mask.

To avoid revisiting regions which have been recently analysed, an inhibition mask can be applied to the last visited locations, both locally for the image and globally for the camera angle: after checking a given area, it is masked so it can not be revisited as soon as the analysis process finishes.

Figure 1 shows an example of this attention model. The top image is the original image with a green rectangle around the most conspicuous region after applying the saliency map. The bottom image is the result of convolving the original image with the saliency map, which darkens the less interesting areas while leaving unchanged the most interesting ones.



Figure 1: Saliency map model

4 Implementation

4.1 Fovea

The fovea is a small depression in the retina (see Fig. 2). It grants the maximum resolution of the whole field of view and despite being

only 1% of the retinal area, more than 50% of visual information processed comes from it. In the rest of the retina, the resolution is inversely proportional to the distance to the fovea. It allows sharp vision and tasks dependant on it such as lecture or driving. The term fovea is often used in Robotics to specify that the center of the image is not treated the same way as the periphery.

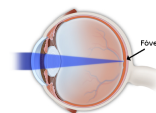


Figure 2: Human eye fovea

Multi resolution sensors (as the fovea, be it real or virtual) are not necessary for the attention system to work. However, there are a couple of reasons (biological and computational) which make interesting to use them:

- The amount of sensor information is reduced in comparison to using the whole field of view at maximum resolution.
- High resolution and wide field of view can be combined as a consequence of the latter.
- Peripheral vision gives only contextual information, allowing for a not so exhaustive process in these areas.

All these advantages can be applied to a robotic vision system, from a hardware or software approach [18, 20]. Hardware solutions are not suitable for us due to our robot being standard, so an algorithmic solution has been chosen.

Our solution is based on C. J. Westelius work [21], but instead of modifying the original image to create the foveal effect, we use a mask to grant access only to certain of its pixels. This mask allows us to simulate the lower resolution of peripheral areas without slowing down the system with unnecessary filter and subsampling operations.

It is possible to configure the amount and size of the multi resolution areas. By default,

we work with three areas: fovea (full resolution, 40% of the image), parafovea (1:2 reduction, 40% of the image) and perifovea (1:4 reduction, 20% of the image). Figure 3 shows the result (right) of applying the mask (center) to the original image (left).

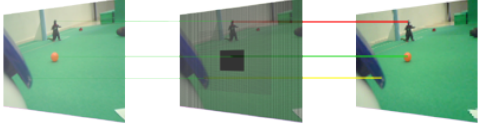


Figure 3: Fovea mask



Figure 4: Original input image

4.2 Maps

The input of the model are 640x480 pixels static RGB color images (Fig. 4 shows the input image which will be used as an example during the explanation of the map construction process). After the previously explained fovea mask has been applied, these images are used to build multiscale pyramids [1] for every map used in the model. Each pyramid has 9 levels and a resolution reduction factor of $1 : 2^n$ for each of them. Level 0 means then no reduction (1:1, original image), while maximum reduction happens at level 8 (1:256). The specific image resolution for each level is then the following:

- level 0: 640x480
- level 1: 320x240
- level 2: 160x120
- level 3: 80x60
- level 4: 40x30
- level 5: 20x15
- level 6: 10x8
- level 7: 5x4
- level 8: 3x2

Intensity maps

The first step of the model consists of creating a nine level intensity pyramid which represents the “intensity” (luminosity) of each image pixel. Using the original image, a intensity matrix M_I is obtained by combination of the R, G and B channels value:

$$m_I(i, j) = \frac{m_R(i, j) + m_G(i, j) + m_B(i, j)}{3}$$

The intensity pyramid is then created using M_I , with $M_I(n)$ being the intensity matrix corresponding to the n th level of the pyramid. Using the pyramid, six intensity maps are obtained by across-scale difference, \ominus , which is obtained by interpolation of the maps to the finer scale and point-by-point subtraction:

$$M_{I(2,5)} = |M_{I(2)} \ominus M_{I(5)}|$$

$$M_{I(2,6)} = |M_{I(2)} \ominus M_{I(6)}|$$

$$M_{I(3,6)} = |M_{I(3)} \ominus M_{I(6)}|$$

$$M_{I(3,7)} = |M_{I(3)} \ominus M_{I(7)}|$$

$$M_{I(4,7)} = |M_{I(4)} \ominus M_{I(7)}|$$

$$M_{I(4,8)} = |M_{I(4)} \ominus M_{I(8)}|$$

This across-scale difference between maps allows for detecting locations at center (areas at scale 2,3,4) which stand out from their surround (scale 5,6,7,8), the same way it happens in human retina [10]. Using several scales for center and surround, instead of just one for each of them, yields truly multiscale feature

extraction [10]. The finest scale is $n = 2$ and not $n = 0$ to reduce noise, excessive detail, and the amount of pixels to be computed (160x120 at scale 2 instead of 640x480 at scale 0), improving both performance and robustness.

Finally, the intensity map I , representing those conspicuous locations from an intensity point of view, is generated combining all the previous maps through across-scale addition, \oplus , which consists of reduction of each map to scale $n = 4$ (40x30 resolution) and point-by-point addition:

$$I = \oplus M_{I(m,n)}$$

Figure 5 shows the intensity map I for the image at Fig. 4.

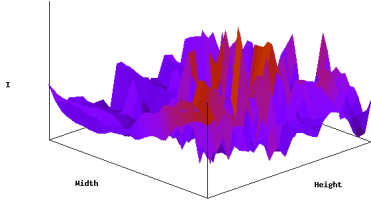


Figure 5: Intensity map

Color maps

Four pyramids representing “color” of each image pixel are created using the normalised R, G and B channels and a yellow channel Y (obtained using the three previous ones): RGB color space channels include intensity information, thus, in order to make the result independent to environmental light, they have to be normalised by intensity. To do so, we applied the same formulae used in [8]

The four color pyramids are used to generate a set of 12 color maps, six for difference between red and green components, $M_{RG(m,n)}$, and six for blue and yellow difference, $M_{BY(m,n)}$, in a similar fashion to the intensity maps.

$$M_{RG(2,5)} = |(M_{R(2)} - M_{G(2)}) \ominus (M_{R(5)} - M_{G(5)})|$$

$$M_{RG(2,6)} = |(M_{R(2)} - M_{G(2)}) \ominus (M_{R(6)} - M_{G(6)})|$$

$$M_{RG(3,6)} = |(M_{R(3)} - M_{G(3)}) \ominus (M_{R(6)} - M_{G(6)})|$$

$$M_{RG(3,7)} = |(M_{R(3)} - M_{G(3)}) \ominus (M_{R(7)} - M_{G(7)})|$$

$$M_{RG(4,7)} = |(M_{R(4)} - M_{G(4)}) \ominus (M_{R(7)} - M_{G(7)})|$$

$$M_{RG(4,8)} = |(M_{R(4)} - M_{G(4)}) \ominus (M_{R(8)} - M_{G(8)})|$$

$M_{BY(m,n)}$ are obtained in a similar way to $M_{RG(m,n)}$ but using the Blue and Yellow components instead.

Finally, a color map C , representing those conspicuous locations from a color contrast point of view, is generated combining all the previous maps:

$$C = \oplus [RG_{I(m,n)} + BY_{I(m,n)}]$$

Figure 6 shows the color map C for the image at Fig. 4.

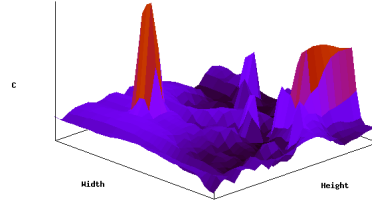


Figure 6: Color map

Orientation Maps

The original model builds up a set of orientation maps which are merged in a final orientation map O which represents the location of those elements which stand out from an orientation point of view in comparison to the rest of the objects present in the image.

Such maps have not yet been implemented in the current version, mainly due to the fact that they are not so important for a controlled

environment like ours (Robocup SPL) in which colour and intensity are already very conspicuous by themselves.

Normalisation

Before obtaining the final saliency maps, all maps have to be normalised.

The Color and Intensity maps obtained are normalised to the same static range $[0..M]$ in order to compare them. Modality dependant differences would also have to be removed. However, since we do not compute orientation maps, this step is not necessary: a 5% intensity difference between two pixels can not be a priori compared to a 0.2 rad orientation difference, but color and intensity differences can be compared without further modification.

A mechanism to promote maps with a small number of strong peaks of activity (conspicuous locations) is also applied. It consists of finding the map's global maximum (M) and computing the average of all its other local maxima (m), globally multiplying the map by $(M - m)^2$. The biggest advantage of this method is its simplicity and speed, while the major drawback is that if a map has two important locations it will only promote the most conspicuous one, hiding the other (humans would probably attend to both of them instead).

In [10], a more complex and efficient method for normalisation based on DoG (Difference of Gaussians) filters is proposed, but it has not yet been implemented.

Saliency map

Once the color and intensity maps have been obtained and normalised, they are combined in the final saliency map S which will guide attention to the most relevant location in the field of view:

$$S = \frac{I+C}{2}$$

Fig. 7 shows the 3D (left) and 2D (right) saliency map S for the image at Fig. 4.

The saliency map S is then applied to the original image as obtained by the robot camera, promoting the most relevant locations and hiding the rest. In Fig. 8 this process is illustrated: left image is the original coloured image. Central image shows the results of apply-

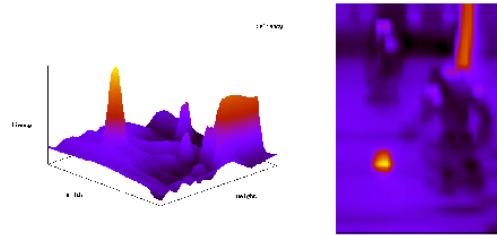


Figure 7: Saliency map

ing the saliency map in Fig. 7 to the original image (the darker the area, the less salient it is). Right image shows the regions with higher saliency across the whole map (green rectangles). Please note that the system proposed only tell us “where” to look at (area) and not “what” (object) to look for; the fact that the ball and the keep are in those areas is a consequence of being the most notorious regions of the image from a color and intensity point of view.



Figure 8: Saliency map applied to the original image

5 Experiments

To test the effectiveness of our approach, the “any ball” Robocup challenge has been chosen. For this challenge, the robot is placed in the game field along with a couple of random coloured and multi-sized balls. The robot has then a couple of minutes to score the biggest amount of goals possible. Classic color filter algorithms used for image segmentation are not usefull in this scenario, since not only ball

color is unknown, but they can also have the same color as the ground (green).

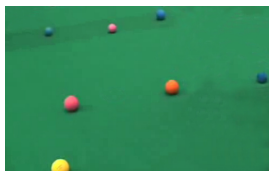


Figure 9: Any ball challenge input image

To simulate this scenario for our system we have given the robot some pictures of the game field containing a random number of different color balls (see Fig. 9).

The model proposed always finds the most salient region in the image, and as long as that region is not dealt with (or inhibited), it will not find any other region. This means that regions chosen as most salient which do not contain any ball must be masked (inhibited), so that others containing a ball can be chosen as focus.

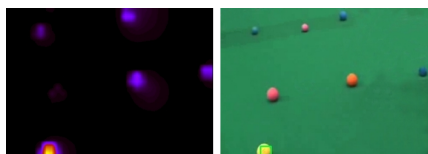


Figure 10: Most salient region of the input image

Once a region containing a ball is chosen (see Fig. 10), robot should approach to it and try to score a goal by kicking it. This part of the experiment has not been implemented yet, but it can be assumed that the ball will end up further from the robot than it was when chosen as focus. To simulate that, once a region containing a ball is chosen by the model, it is assumed that the robot could kick it and that specific ball is removed from the next input image for the robot.

With the originally most salient ball no longer present in the field of view (see Fig. 11), the saliency map changes and a new most salient region is chosen (see Fig. 12). The previously explained process is now repeated: if

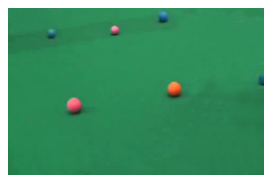


Figure 11: Any ball challenge second input image

the new region contains a ball, it is chosen as focus and kicked, otherwise, it is inhibited and the second most salient region is checked, repeating the process until finding a region containing a ball or not finding any at all.

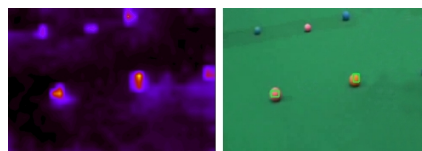


Figure 12: Most salient regions of the second input image

The results obtained are very promising, with a 100% success rate for the images used. Even the regions containing small balls with almost the same color of the ground are chosen in the last iterations of the algorithm (see Fig. 13). It can be easily understood that the color map (top left image at Fig. 14) gives no useful information in this case, since the whole field of view is almost of the same color (except for the lines). However, the intensity map (top right image at Fig. 14) shows strong peaks at those areas containing either shades, which should be minimal except for the one belonging to the ball (due to it being the only object in the field apart from the robot), or different light reflection patterns, as it happens with the region containing the ball since the ball is made of a different material from the ground's. The final saliency map obtained once again chooses the region containing the ball as the most salient one (see bottom left and bottom right images at Fig. 14).



Figure 13: Input image containing a ball of the same color of the ground

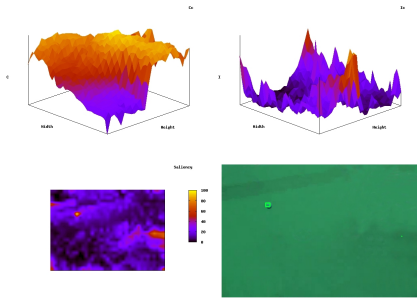


Figure 14: Color, intensity and saliency map and most salient regions of an input image containing a ball of the same color of the ground

6 Discussion and further work

In this paper we have presented an attention control model based on a saliency map which mainly differs from the original by Itti [10] in two aspects: the saliency map is obtained using only intensity and color information, dispensing orientation and movement data, and input images' size is reduced by using a fovea mask. These modifications improve the model performance and allow for a better adaptation to the Robocup SPL environment.

The model has proved to be useful for the “any ball” challenge, with better results than classic filter and segmentation algorithms, which do not provide results robust enough when trying to identify balls of similar color to the field.

The main drawback of our proposal is the time it consumes, which makes the model not usable for real time game play. However, the system remains suitable for competition when combined with classic color filter algo-

rithms, applying the saliency calculation only to certain images or situations (finding areas in the field containing interesting objects, for instance a ball in the proposed challenge) and using the classic color filter approach for the rest of the tasks (object recognition and subsequent tracking).

There are mainly two topics which would need to be addressed in the near future: a more effective normalisation operator and time consumption optimisation.

The simple normalisation operator used tends to promote only one activity peak in the intermediate maps, which makes the most conspicuous area hide the rest even if there is a second one very close to it (and thus also very important from a saliency point of view). This leads to occasional problems. For instance, when both the ball and the yellow keep are visible, specially with partial ball occlusions, the yellow net may hide the ball in the final saliency map. In the “any ball” challenge experiment here explained, it can be seen that, for the same reason, some of the regions containing balls are not found until second iteration (compare Fig. 10 to Fig. 12) when the previously most salient region (the one containing the yellow ball) has been removed. Itti et al. already solved this issue by using DoG filters instead [10], which makes the system work better in these cases.

As previously stated, time consumed by the maps generation algorithm is too high. One of the main advantages of attention is the great reduction in the amount of information to process, specially since processing a stream of video in limited hardware as a robot is a high time-consuming task. However, the whole process is taking around 200 ms, which is an excessive amount of time to make it worthwhile in this sense. An optimisation of the code could make the system much more suitable for full time use.

Acknowledgment

The authors would like to thank the Spanish Ministry of Innovation for its support to this project under the grant DPI2007-66556-

C03-01 (COCOGROM project), and Junta de Castilla y León and European Social Fund for their support to Juan F. Garcia.

References

- [1] C. Anderson and D. Van Essen. Shifter circuits: a computational strategy for dynamic aspects of visual processing. In *Proc. Nat. Acad. Sci. USA*, 1987.
- [2] J. R. Anderson. *Cognitive psychology and its implications*. Worth Publishers, 2004.
- [3] P. Bachiller, P. Bustos, and L. J. Manso. Attentional selection for action in mobile robots. *Advances in Robotics, Automation and Control*, 2008.
- [4] D. E. Broadbent. *Perception and communication*. Pergamos Press, New York, 1958.
- [5] G. Deco. *A Neurodynamical Model of Visual Attention: Feedback Enhancement of Spatial Resolution in a Hierarchical System*. 2000.
- [6] S. Frintrop, G. Backer, and E. Rome. Goal-directed search with a top-down modulated computational attention system. *Lecture Notes in Computer Science*, 3663:117–124, 2005.
- [7] Jacob M. Gryn, Richard P. Wildes, and John K. Tsotsos. Detecting motion patterns via direction maps with application to surveillance. *Computer Vision and Image Understanding*, 113:291–307, 2009.
- [8] L. Itti and C. Koch. A saliency-based research mechanism for overt and covert shifts of visual attention. *Vision Research*, 40:1489–1506, 2000.
- [9] L. Itti and C. Koch. Computational modelling of visual attention. *Nature Reviews Neuroscience*, 2:194–203, 2001.
- [10] L. Itti, C. Koch, and E. Niebur. Attentive mechanisms for dynamic and static scene analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20:1254–1259, 1998.
- [11] L. Itti, C. Koch, and E. Niebur. A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20:1254–1259, 1998.
- [12] B. Julesz. Early vision and focal attention. *Review of Modern physics*, 66(3):735–772, 1991.
- [13] C. Koch and S. Ullman. Shifts in selective visual attention: Towards the underlying neural circuitry. *Human Neurobiology*, 4:219–227, 1985.
- [14] D. LaBerge, R.L. Carlson, J. K. Williams, and B. G. Bunney. Shifting attention in visual space: Tests of moving-spotlight models versus an activity-distribution model. *J. Experimental Psychology: Human Perception and Performance*, 23:1380–1392, 1997.
- [15] O. Neumann, A. H. C. van der Heijden, and A. Allport. Visual selective attention: introductory remarks. *Psychological Research*, 48:185–188, 1986.
- [16] B. Olshausen, C. Anderson, and C.V. Essen. A multiscale dynamic routing circuit for forming size- and position-invariant object representations. *J. Computational Neuroscience*, 2:45–62, 1995.
- [17] H. Pashler and P. Badgio. Visual attention and stimulus identification. *J. Experimental Psychology*, 11:105–121, 1985.
- [18] M. Tistarelli and G. Sandini. Direct estimation of time-to-impact from optical flow. *IEEE Workshop on Visual Motion*, pages 52–60, 1991.
- [19] A. Torralba, A. Oliva, M. S. Castellanos, and J. M. Henderson. Contextual guidance of eyes movements and attention in real-world scenes: the role of the global features in object research. *Psychological Review*, 113:766–786, 2006.
- [20] J. van der Spiegel, G. Kreider, C. Claeys, I. Debusschere, G. Sandini, P. Dario, F. Fantini, P. Bellutti, and G. Soncini. *A foveated retina like sensor using ccd technology*. Kluwe, 1989.
- [21] C.J. Westelius. *Focus of attention and gaze control for robot vision*. PhD thesis, Department of Electrical Engineering, Linköping University, Sweden, 1995.
- [22] J. M. Wolfe. *Guided Search 4.0: Current Progress with a model of visual search*. Brigham and Womens Hospital and Harvard Medical School, 2007.
- [23] R.D. Wright and L.M. Ward. *Orienting of Attention*. Oxford University Press, 2008.