



PERGAMON

Vision Research xxx (2002) xxx-xxx

Vision  
Research

www.elsevier.com/locate/visres

## Eye movements in iconic visual search

Rajesh P.N. Rao<sup>a</sup>, Gregory J. Zelinsky<sup>b</sup>, Mary M. Hayhoe<sup>c</sup>, Dana H. Ballard<sup>d,\*</sup>

<sup>a</sup> Department of Computer Science, University of Rochester, Rochester, NY 14627, USA

<sup>b</sup> Beckman Institute, University of Illinois, Urbana, IL 61801, USA

<sup>c</sup> Center for Visual Science, University of Rochester, Rochester, NY 14627, USA

<sup>d</sup> Department of Computer Science, University of Rochester, Rochester, NY 14627, USA

Received 15 January 2001; received in revised form 17 July 2001

### Abstract

Visual cognition depends critically on the moment-to-moment orientation of gaze. To change the gaze to a new location in space, that location must be computed and used by the oculomotor system. One of the most common sources of information for this computation is the visual appearance of an object. A crucial question is: How is the appearance information contained in the photometric array is converted into a target position? This paper proposes a model that accomplishes this calculation. The model uses iconic scene representations derived from oriented spatiochromatic filters at multiple scales. Visual search for a target object proceeds in a coarse-to-fine fashion with the target's largest scale filter responses being compared first. Task-relevant target locations are represented as saliency maps which are used to program eye movements. A central feature of the model is that it separates the targeting process, which changes gaze, from the decision process, which extracts information at or near the new gaze point to guide behavior. The model provides a detailed explanation for center-of-gravity saccades that have been observed in many previous experiments. In addition, the model's targeting performance has been compared with the eye movements of human subjects under identical conditions in natural visual search tasks. The results show good agreement both quantitatively (the search paths are strikingly similar) and qualitatively (the fixations of false targets are comparable). © 2002 Published by Elsevier Science Ltd.

**Keywords:** Saccades; Computation; Attention; Visuomotor control

### 1. Introduction

Human vision relies extensively on the ability to make saccadic eye movements to orient the high-acuity foveal region of the eye over targets of interest in a visual scene. However, resolution per se is not the only determinant of gaze location. Starting from Yarbus' classical work (Yarbus, 1967), many studies have suggested that gaze changes are directed according to the ongoing cognitive demands of the task at hand. The task-specific use of gaze is best understood for reading text (O'Regan, 1990) where the eyes fixate almost every word, sometimes skipping over small function words. In addition, saccade size during reading is modulated according to the specific nature of the pattern recognition task at hand (Kowler & Anton, 1987). Tasks requiring comparison of complex patterns also elicit characteristic saccades back

and forth between the patterns (Just & Carpenter, 1976). In copying of a model block pattern on a board, subjects have been shown to employ fixations for accessing crucial information during different stages of the task (Ballard, Hayhoe, & Pelz, 1995; Ballard, Hayhoe, Pook, & Rao, 1997). In natural language processing, fixations can reflect the instantaneous parsing of a spoken sentence in the current visual context (Tanenhaus, Spivey-Knowlton, Eberhard, & Sedivy, 1995). The role of gaze has been studied by in a variety of natural visuomotor tasks such as driving, music reading and playing ping-pong (Land & Furneaux, 1997). In each case, gaze was found to play a central *functional* role, closely linked to the immediate task demands. All these tasks have very different kinds of fixation targets, sometimes only defined in terms of functional needs. For example, in driving around a bend, subjects fixate the tangent point of the curve to control steering angle, and in ping-pong, subjects fixate the bounce point in advance, in order to estimate the ball's trajectory.

\* Corresponding author. Tel.: +1-716-275-3772; fax: +1-716-461-2018.

E-mail address: dana@cs.rochester.edu (D.H. Ballard).

59 The general utility of saccadic eye movements has  
60 spurred an extensive effort to characterize their proper-  
61 ties. A variety of studies have revealed the importance of  
62 task, acuity, and visual features in determining the  
63 stimulus for target selection together with accompany-  
64 ing metrics of accuracy and fixation duration (e.g.  
65 Findlay, 1997; Hooge & Erkelens, 1998; Motter &  
66 Belky, 1998; Viviani, 1990; Zelinsky & Sheinberg, 1997).  
67 However, much less is known about the underlying  
68 computational operations that determine these proper-  
69 ties, although some ground-breaking work has been  
70 done. Itti and Koch (2000) use the coincidental align-  
71 ment of visual features to define a saliency map of  
72 possible targets. Moving the gaze to these points suc-  
73 cessively has some resemblance to human visual search  
74 but there is no model of how specific targets are selected.  
75 Tsotsos et al. (1995) use an hierarchical attractor net-  
76 work to define interesting targets. Unlike Itti and Koch,  
77 Tsotsos's network can be driven by selected target fea-  
78 tures, however the representation cannot define com-  
79 pletely general image targets. There also has been no  
80 attempt in either of these models to compare their per-  
81 formance with human visual search.

82 This paper describes a general model for fixating and  
83 remembering appearance-based encodings of targets in  
84 natural scenes. The model uses iconic (appearance-  
85 based) target representations to search arbitrary visual  
86 scenes. Iconic representations are specified by the re-  
87 sponses of oriented spatiochromatic filters at multiple  
88 scales. This has been demonstrated to be a very robust  
89 computational mechanism for target selection in natural  
90 scenes (Rao & Ballard, 1997). The computation of target  
91 coordinates for a saccade reduces to correlation between  
92 a "top-down" iconic target representation and the  
93 "bottom-up" iconic scene representations. The model  
94 provides a good fit to visual search data where the target  
95 is defined predominantly from its appearance. An key  
96 feature of the model is that it separates the targeting  
97 process, which changes gaze, from the decision process,  
98 which uses the information at the new gaze point. The  
99 virtue of this separation is that decision-making about  
100 the target can be separated from the process of fixating  
101 it. Thus there is no additional control structure to make  
102 the gaze change contingent on the decision process. If  
103 the decision process is slow with respect to the time  
104 needed for target selection, then gaze can be moved to  
105 the target more accurately. If the decision process is fast,  
106 then gaze does not have to be changed at all, as is ob-  
107 served in a huge number of studies of attention.

## 108 2. General purpose iconic representations

109 In many experiments that study saccades, the targets  
110 themselves are simple colored shapes that are presented  
111 on a blank background. While extensive useful data has

been collected using this paradigm, this setup does not  
address issues of target selection in natural viewing. In  
natural scenes, the saccadic target may be composed of  
complex photometric intensity patterns, produced by  
cluttered scenes. In order to move the eyes in this case,  
there must be a mechanism that translates the intensity  
image on the retina into a representation that can be  
used by the oculomotor system. Such a mechanism must  
meet at least the following three criteria:

1. *Generality*: Any proposed mechanism for targeting parts of an image must have broad generality since saccadic targets can vary greatly according to the requirements of the current task.
2. *Speed*: Targets must be computed quickly in order to model observed human performance. Using millisecond neural circuitry, the targets for the next fixation need to be computed in approximately 80–100 ms, allowing barely one pass through the cortex (Oram & Perrett, 1992; Thorpe & Imbert, 1989).
3. *Resolution*: The computation of the target must use spatial scales that are available extrafoveally, since it is unlikely that the target is already at the gaze point.

One representation that meets these criteria employs low resolution iconic representations of targets and scenes that can be extracted directly from the optic array. This allows general portions of a scene to be represented in a precategorical format without requiring any elaborate segmentation. This is an essential property, since the information required for such complex operations is frequently the goal of the eye movement itself. The computation of saccadic target coordinates is accomplished by correlating the iconic memory of the target with the iconic representation of the current optic array. A correlation peak indicates the most likely location of the target in the current image, allowing a saccade to be executed to that location. We regard the notion of "icon" as completely general. The idea is that any criterion for a gaze point can be transformed into an appearance model which captures how that criterion should appear in the scene. Then the resultant appearance image, or icon, is used as a correlation template.

It would be prohibitively expensive to encode icons literally as gray-level images, since the memory needed would then scale with the size and number of icons. A more efficient alternative is to encode the icons as their responses to a set of spatiochromatic basis functions, or spatial filters (Itti & Koch, 2000; Poetzsch, Krueger, & Von der Malsburg, 1996; Weber & Malik, 1995). One motivation for this is that it approximates the transformations imposed by the receptive fields of striate cortical cells. Another motivation is the psychophysical evidence of suggesting that the human visual system uses such channels (Graham, 1989; Wilson & Wilkinson,

1997). The particular filters we use are the steerable filters, so-called because the responses of these filters at any given orientation can be used to produce the responses at any other location by interpolation formulae. A local image patch can be characterized using a zeroth order Gaussian  $G_0$  and nine of its oriented derivatives (Fig. 1) as follows (Freeman & Adelson, 1991):

$$G_n^{\theta_n}, \quad n = 1, 2, 3, \quad \theta_n = 0, \dots, m\pi/(n + 1),$$

$$m = 1, \dots, n \quad (1)$$

where  $n$  denotes the order of the filter and  $\theta_n$  refers to the preferred orientation of the filter. The response of an image patch  $I$  centered at  $(x_0, y_0)$  to a particular basis filter  $G_i^{\theta_j}$  can be obtained by convolving the image patch with the filter:

$$r_{i,j}(x_0, y_0) = \int \int G_i^{\theta_j}(x_0 - x, y_0 - y) I(x, y) dx dy \quad (2)$$

The iconic representation for the local image patch centered at  $(x_0, y_0)$  is formed by combining into a high-dimensional vector the responses from the 10 basis filters above at different scales

$$\mathbf{r}(x_0, y_0) = [r_{i,j,s}(x_0, y_0)] \quad (3)$$

where  $i = 0, 1, 2, 3$  denotes the order of the filter,  $j = 1, \dots, i + 1$  denotes the different filters per order, and  $s = s_{\min}, \dots, s_{\max}$  denotes the different scales of the filters. For computational efficiency, a Gaussian pyramid representation of the image can also be used to generate multi-scale responses from a set of basis filter kernels at a fixed scale. This strategy was used in the visual search simulations. As an example, Fig. 2 shows the filter-based responses at a given location in a cluttered scene for filters  $G_1$  and  $G_2$  and five spatial scales. The filter response vector at every image location in

general provides an almost unique representation of the local image region surrounding that location (Rao & Ballard, 1996).

The model search simulations used gray scale stimuli, with three spatial scales and nine filters per scale for a total of 27 measurements per image location. The scales used in our tests range from approximately 1–6 cycles per degree, well within the limits of human spatial resolution at the eccentricities involved in the experiments described here. The basis functions described above were picked a priori, but very similar functions can be learned from samples of natural images (Ballard et al., 1997; Barrow, 1987; Bell & Sejnowski, 1997; Hancock et al., 1992; Olshausen & Field, 1996).

The use of multiple scales is crucial to the visual search model. In particular, the larger the number of scales, the greater the perspicuity of the representation as depicted in Fig. 3, which shows the frequency distribution of correlations between all points in the dining table image (Fig. 8(d)) and a fixed target point in the same image. The distribution on the left shows how using filter responses at a single scale causes ambiguity in the iconic scene representations, with as many as 936 points in the scene having correlations greater than 0.94 with respect to a fixed target. However, when five scales are used, the ambiguity is resolved, and only a single point (the target point) correlates greater than 0.94 (indicated by the arrow for both histograms). The greater perspicuity results partly due to the inclusion of information from additional scales and partly due to the high-dimensionality of the multi-scale vectors. The high-dimensionality of the vectors makes them remarkably robust to noise due to the *orthogonality* inherent in high-dimensional spaces: given any vector, almost all of the other vectors in the space tend to be relatively uncorrelated with the given vector (Kanerva, 1988; Rao & Ballard, 1995a), and almost none are identical with respect to each other. The result is that the filter response vector for a given point is unique for all practical purposes and can therefore be used to define search targets. This property also makes the filter template robust to partial occlusions, which commonly occur in natural viewing (see Rao & Ballard (1995a) for some examples).

The representation works best when the gross viewpoint of the scene does not change drastically from moment-to-moment. The filter responses are dominated by a cosine envelope, so that there is a useful range of rotations for which the responses will be effectively invariant. Drastic rotations are handled by storing feature vectors from different views (Bulthoff & Edelman, 1992). This is consistent with psychophysical evidence that shows that subjects represent objects using a small number of separate viewpoints. The multi-scale representation also allows interpolation strategies for scale invariance (Rao & Ballard, 1995a).

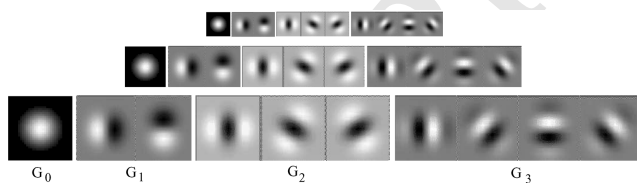


Fig. 1. Spatiochromatic basis functions. Motivation for these basis functions comes from statistical characterizations of natural image stimuli (Bell & Sejnowski, 1997; Derrico & Buchsbaum, 1991; Hancock, Baddeley, & Smith, 1992; Olshausen & Field, 1996; Rao & Ballard, 1997). The nine oriented spatial filters at three octave-separated scales for each of the three channels in (a) (bright regions denote positive magnitude while darker regions denote negative magnitude). At each scale, these nine filters are comprised of two first-order derivatives ( $G_1$ ) of a 2D photometric Gaussian, three second-order derivatives ( $G_2$ ), and four third-order derivatives ( $G_3$ ). Thus, there are three scales per channel, and nine spatial filters per scale, for a total of 27 filter responses characterizing each location in the image. These 27 spatiochromatic measurements at a given image location can be regarded as a photometric signature of the local image region centered at that location.

196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250

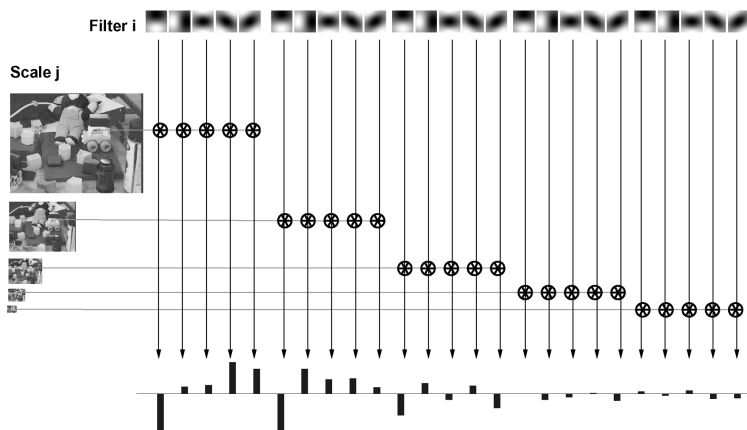


Fig. 2. Using spatiochromatic filters to extract task-dependent properties. A portion of a cluttered image. The scales at which the filters of Fig. 1 were applied to the image are shown on the left. Each individual filter, when convolved with the local image intensities near the given image location, results in one measurement. This example uses the first two filters and five spatial scales for a total of 25 measurements per point. Positive responses in the vector are represented as an upward bar above the horizontal, negative responses as a downward bar below the horizontal. For reasons of economy, large scale filters are modeled by using the standard size filter and shrinking the image.

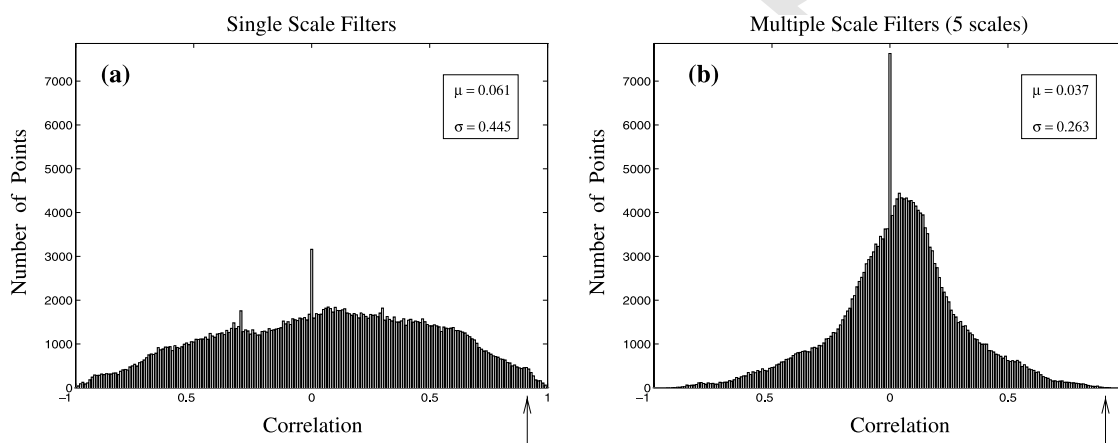


Fig. 3. The effect of scale. The distribution of distances (in terms of correlations) between the filter response vector for a selected target point in the dining table scene (Fig. 7(a)) and all other points in the scene is shown for single scale response vectors (a) and multiple scale vectors (b). Using responses from multiple scales (five in this case) results in greater perspicuity and a sharper peak near 0.0. The most important feature of these plots appears at the extreme right hand side. Only one point (the target point) has a correlation greater than 0.94 (demarcated by an arrow) in the multiple scale case (b) whereas 936 candidate points fall in this category in the single scale case (a).

251 To summarize, the representation meets the criterion  
 252 of generality since any gaze target can be translated in-  
 253 ternally into a local appearance, which in turn can be  
 254 expressed in terms of filter responses. The representation  
 255 can be used quickly since targeting reduces to filter  
 256 correlations, which we assume can be done in parallel  
 257 without penalty over the retinal array. Finally the use of  
 258 multiple scales means that the range of resolutions used  
 259 can be adjusted to trade-off speed with accuracy as  
 260 suggested by Geisler and Chou (1995).

### 3. Modeling visual search

261

262 Early models of visual search suggested that the  
 263 search process proceeds item-by-item (Treisman, 1988)  
 264 but data showing fast search times for some multiple  
 265 conjunctions were hard to model. More recent models,  
 266 guided by Palmer, Vergese, and Pavel (2000) assume  
 267 that search is area-based, aimed at detecting targets  
 268 within a window centered around the center of gaze  
 269 (Eckstein, 1998; Geisler & Chou, 1995). The size of the  
 270 window is a function of the speed and accuracy required  
 271 of the task, and reflects the signal-to-noise characteris-  
 272 tics of the display (Motter & Belky, 1998). In the latter

273 case, the search task can be seen as one of covering the  
274 scene while prioritizing likely locations. As a conse-  
275 quence the gaze point need not search item-by-item but  
276 can delimit large areas.

277 Fig. 4 motivates the model's use of area-based search  
278 in terms of the resolution of the retinal image as re-  
279 ported by Hess. For each search task, a resolution needs  
280 to be chosen based on signal-to-noise conditions of the  
281 display and the spatial properties of the target. The  
282 resolution chosen for the search process defines a search  
283 window width. Higher signal-to-noise means that the  
284 object can be recognized at a lower resolution and hence  
285 a bigger search window can be used. A consequence of  
286 this choice is that the same resolution is used throughout  
287 the search window, even though higher resolution is  
288 available. The use of a set resolution in this manner by  
289 our model is counterintuitive, as it is more natural to  
290 assume that all the available resolution is continuously  
291 available. However, the use of resolution as a search  
292 parameter is motivated by search experiments that show  
293 that other search parameters are set and changed with  
294 temporal cost. For example, Sperling (Sperling & Do-  
295 sher, 1986) showed that searching displays of two dif-  
296 ferent font sizes incurred a cost that suggested the scale  
297 had to be set for each size.

298 The visual search model is composed of three separate  
299 procedures that each operate largely independent of  
300 each other, while at the same time cooperating to solve  
301 the current visual search task:

1. A *targeting process* (or “where” process) that com- 302  
putes the next location to be fixated. 303
2. A *decision process* (or “what” process) that matches a 304  
stored iconic object representation to the current fov- 305  
eated image region. 306
3. An *oculomotor process* that accepts retinotopic target 307  
locations from the “where” process and executes a 308  
saccade to the target location (a method for learning 309  
this sensorimotor mapping is given in (Rao & Bal- 310  
lard, 1995b)). 311

312 The model assumes that these processes are running 312  
concurrently, but that they do not have to be precisely 313  
coordinated in time. The oculomotor process will con- 314  
tinue to execute eye movements as long as the decision 315  
process has not terminated. The current best guess of 316  
target location is updated as fixations increase the 317  
available resolution. Although we do not model the 318  
decision process, a key point is that the decision process 319  
needs to choose a resolution and window in the same 320  
way as the search process, but the two resolutions need 321  
not be the same, since getting the gaze to the target and 322  
analyzing a property of the target are different compu- 323  
tations. 324

325 All three processes use a *saliency map* (Koch & Ull- 325  
man, 1985) whose value at a given location represents 326  
the weight determined by multi-scale filter-based corre- 327  
lation. This weight map has a dual purpose: (1) it allows 328  
the oculomotor process to fixate target locations with 329  
high correlations, and (2) its maximum value is used by 330

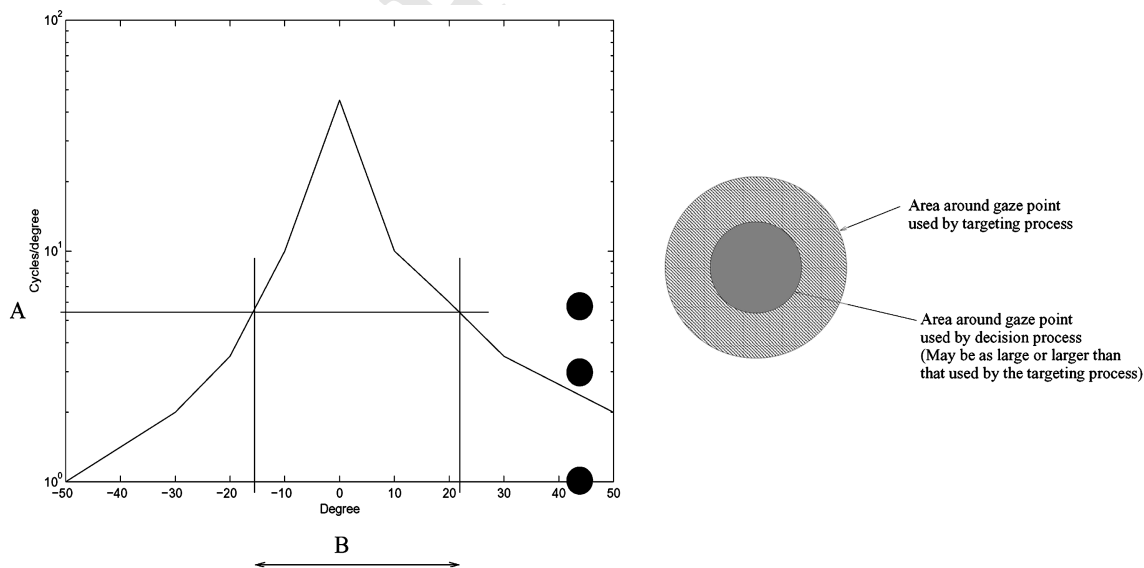


Fig. 4. How the model chooses resolutions. Left: Resolution as a function of retinal eccentricity, with a hypothetical search window. Data are replotted from (Anderson, Mullen, & Hess, 1991). For a given search task our model assumes that the subject chooses a signal-to-noise ratio. That defines a maximum resolution to be used in the search (A). Given this resolution value, the resolution available on the retina defines a search width (B). The three frequency scales used by the model are shown at right as filled circles. Right: Separate search windows are used for targeting, which changes gaze, and decisions, which extract information needed for behavior.

331 the decision process to judge the presence or absence of  
332 the target. The decision process need only use a signal-  
333 to-noise criterion to decide whether the correlation peak  
334 in the saliency map is high enough so that the target can  
335 be assumed to be present. It does not need information  
336 on where that measurement came from.

337 The computation of such a saliency map usefully can  
338 be described in an oversimplified form as follows. Ob-  
339 jects of interest to the current search task are assumed to  
340 be represented by a set of memorized filter response  
341 vectors  $\mathbf{r}_s^m$  where  $s$  denotes the scale of the filters and  $m$   
342 denotes a particular target object in memory. Given a  
343 new input image, the targeting process computes the  
344 most likely location of the target as follows:

- 345 1. Compute the saliency map  $S$  across all locations  $(x, y)$   
346 as

$$S(x, y) = \sum_{s=1}^{\max} \|\mathbf{r}_s(x, y) - \mathbf{r}_s^m\|^2 \quad (4)$$

348 where  $\|\mathbf{x}\|$  denotes the Euclidean norm of the vector  
349  $\mathbf{x}$ . In other words, the saliency value at location  $(x, y)$   
350 is simply the sum of squared differences between the  
351 corresponding components of the filter response  
352 vector  $\mathbf{r}_s$  at that image location and the memorized  
353 target object vector  $\mathbf{r}_s^m$ , across all filter scales  
354  $s = 1, \dots, \max$ .

- 355 2. The location for saccadic targeting is the one that is  
356 most similar to the target, where similarity is given  
357 by Euclidean distance

$$(\hat{x}, \hat{y}) = \arg \min S(x, y) \quad (5)$$

360 In this targeting process, a single saliency map is  
361 calculated across all filter scales for a given image, and

the location  $(\hat{x}, \hat{y})$  to be foveated is chosen to be the one  
with the highest correlation value with respect to the  
memorized target i.e. the one with the least  $S(x, y)$ .  
These computations have been implemented using the  
Datacube MV200 image processor and the Rochester  
dual-camera robot head to perform targeting move-  
ments in real time in natural scenes. The virtue of this  
system is that the Datacube MV200 can compute con-  
volutions at frame rates ( $30 \text{ s}^{-1}$ ) and this allows for  
extensive experimentation. Details of the hardware im-  
plementation are given in (Rao & Ballard, 1995a). Figs.  
5 and 6 illustrates the utility of this simple algorithm in a  
search task. Gaze, as denoted by the cross-hairs, is first  
directed to a given scene location as shown in (a). At  
that point the filter responses are memorized. Next, at  
some point in the course of the rest of the behavior, it  
may be desirable to return to the original location from a  
distal point. The targeting algorithm is used to correlate  
the memorized features with the current retinotopic  
image, resulting in a saliency map as shown in (c). Note  
that the coordinate system of the saliency map can also  
be interpreted in terms of a motor error signal. Thus, the  
saliency peak can be used to drive the oculomotor  
command for returning the eyes to the original target  
without involving complex object properties.

#### 4. Human fixation patterns in appearance-based visual search

Human fixation patterns are more complicated than  
those predicted by the simple search model. In order to  
compare the model's performance with human search  
and targeting behavior we used the data from eye  
movements in a visual search task described in (Zelin-

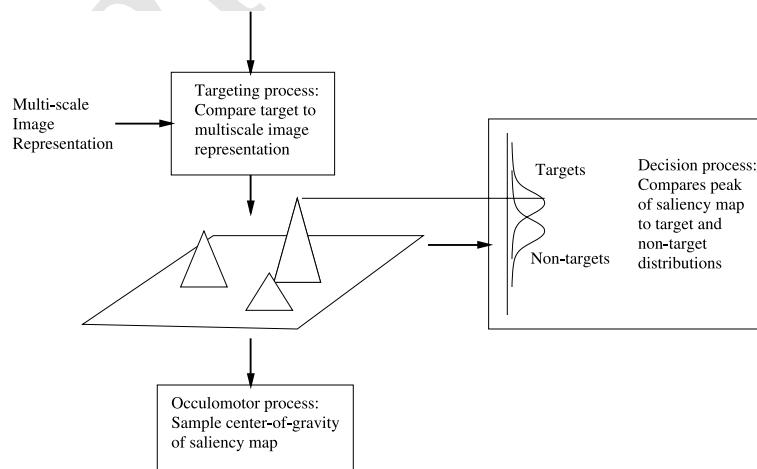


Fig. 5. Visual search using spatial filter responses. The simplest form of the visual search model is based on winner-take-all correlation matching. (a) At a given location, the filter responses are remembered. (b) Next, gaze is transferred to another point. The search problem is to find the original location in this new view. (c) The saliency map, showing the highest correlation value (brightest point) at the original location. (d) Gaze is transferred back to that location.

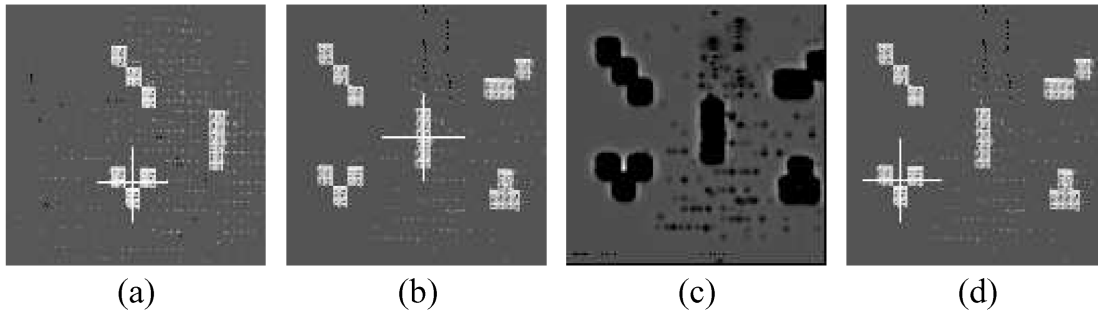


Fig. 6. Visual search using spatial filter responses. The simplest form of the visual search model is based on winner-take-all correlation matching. (a) At a given location, the filter responses are remembered. (b) Next, gaze is transferred to another point. The search problem is to find the original location in this new view. (c) The saliency map, showing the highest correlation value (brightest point) at the original location. (d) Gaze is transferred back to that location.

394 sky, Rao, Hayhoe, & Ballard, 1997). In this experiment,  
395 fixation patterns were observed in a simple search par-  
396 adigm using natural images of three different scenes: a  
397 crib, a workbench and a dining table. Subjects were  
398 asked to fixate a point near the bottom of a  $12^\circ \times 16^\circ$   
399 display. They were given a one second presentation of  
400 an image containing a single object (e.g. a tool) at the  
401 fixation point, defining the search target, on a realistic  
402 background (e.g. the workbench). This was followed  
403 approximately one second later by a scene that filled the  
404 display and contained one, three, or five objects (e.g.  
405 various tools) on the same background. Images of the  
406 objects were placed on the background on-line at one to  
407 five of the six possible equi-eccentric locations ( $22.5^\circ$ ,  
408  $45^\circ$ ,  $67.5^\circ$ ,  $112^\circ$ ,  $135^\circ$ , and  $157.5^\circ$ , each located at an  
409 eccentricity of  $7^\circ$ ) along an arc centered on the subject's  
410 initial fixation point (see Fig. 7(a)). The objects them-  
411 selves subtended about  $2^\circ$  of visual angle. The subjects  
412 were asked to indicate (by pressing a button), as quickly  
413 and accurately as possible, whether the previewed object

was among the group of one to five objects in the sub- 414  
sequent view. Note that the configuration of the objects 415  
in the experiment was like that shown in the following 416  
figure (see Fig. 7(a)). For each subject, each of the search 417  
trials tested a unique configuration of objects and po- 418  
sitions. The trials were evenly divided into randomly 419  
interleaved target-present and target-absent conditions 420  
for set sizes of one, three, and five objects. The back- 421  
ground objects were always present. Eye movements 422  
were recorded when the subjects performed this visual 423

The typical eye movements elicited in this particular 432  
task are shown in Fig. 7(a). The surprising result was 433

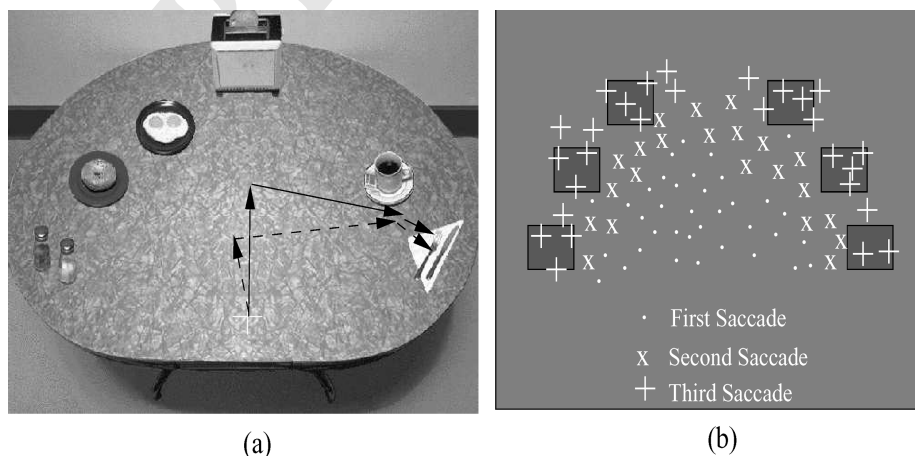


Fig. 7. Eye movements in the visual search task. Measurements from actual human data show marked differences from the simple winner-take-all model: (a) shows the typical pattern of multiple saccades (shown here for two different subjects) elicited during the course of searching for the object composed of the fork and knife. The initial fixation point is denoted by "+"; (b) depicts a summary of such movements over many target-present search trials as a function of the six possible locations of a target object on the table.

434 that rather than a single movement to the location of the  
435 memorized target, several saccades are typical, with each  
436 successive saccade moving closer to the target location  
437 (Fig. 7(b)). This “skipping” of the saccades in this  
438 search paradigm proved to be an extraordinarily robust  
439 finding, occurring in almost all 480 trials across all four  
440 subjects (Zelinsky et al., 1997).

#### 441 5. Appearance-based search model

442 The simple model described in Section 3 cannot ac-  
443 count for the experimentally observed multiple fixations,  
444 since its winner-take-all strategy means that only a sin-  
445 gle saccade is computed. However, multiple fixations  
446 can be fairly easily modeled if the computation of the  
447 saliency map is modified in the following three ways:

448 (1) The saliency map computation is made to be  
449 slower than the time needed to make an eye movement.  
450 This would imply that eye movements are made to tar-  
451 get locations as determined by the *current* state of the  
452 saliency map, rather than waiting until the final state has  
453 been computed.

454 (2) The saliency map is computed using the larger  
455 spatial scale filters first, adding saliency information  
456 from successively finer scales as the search process  
457 evolves over time. Motivation comes both from the data  
458 and several studies that show that lower spatial fre-  
459 quencies influence the decision process earlier than  
460 higher spatial frequencies (Bichot & Schall, 1999; Gil-  
461 christ & Heywood, 1999; McPeck & Keller, 2001;  
462 Schyns & Oliva, 1994).

463 (3) The most likely target location is computed using a  
464 weighted averaging scheme rather than a pure winner-  
465 take-all mechanism. In conjunction with (1) and (2)  
466 above, this would imply that early eye movements are  
467 directed to “center-of-gravity” locations since only  
468 coarse scale information regarding the objects and the  
469 background is available at the early stages of the search,  
470 thereby biasing the weighted averaging model towards  
471 the center of the scene. The motivations for doing this is  
472 that it is known that in some circumstances saccades  
473 display a “center-of-gravity” property and fall midway  
474 between potential targets (Coren & Hoenig, 1972;  
475 Findlay, 1982, 1987; He & Kowler, 1989). The move-  
476 ment of the first saccade to the center of the image is  
477 likely to be a center-of-gravity effect, caused by the  
478 presence of many potential targets in the scene.

479 To implement these modifications, the simple winner-  
480 take-all model of Section 3 was changed to the follow-  
481 ing:

- 482 1. Set the initial scale of analysis  $k$  to the largest scale  
483 i.e.  $k = \max$ ; set  $S(x, y) = 0$  for all  $(x, y)$ .

2. Compute the current saliency map across all locations  
 $(x, y)$  based on filter responses from the current scale  
 $k$  up to the maximum scale

$$S(x, y) = \sum_{s=k}^{\max} \|\mathbf{r}_s(x, y) - \mathbf{r}_s^m\|^2 \quad (6)$$

As before,  $S(x, y)$  is the square of the Euclidean dis-  
tance between the filter response vector  $\mathbf{r}_s$  for image  
location  $(x, y)$  and the memorized target response  
vector  $\mathbf{r}_s^m$ , summed over the scales  $s = k, \dots, \max$ .

3. Find the location for saccadic targeting using the fol-  
lowing *weighted population averaging scheme*:

$$(\hat{x}, \hat{y}) = \sum_{(x,y)} F(S(x, y))(x, y) \quad (7)$$

where  $F$  is an interpolation function. For the experi-  
ments, we used

$$F(S(x, y)) = \frac{\exp(-S(x, y)/\lambda(k))}{\sum_{(x,y)} \exp(-S(x, y)/\lambda(k))} \quad (8)$$

This choice is attractive since it allows an interpre-  
tation of the search algorithm as computing *maxi-  
mum likelihood estimates* (cf. Nowlan, 1990) of target  
locations. In the above,  $\lambda(k)$  is a “temperature” pa-  
rameter that is decreased with  $k$ . Decreasing  $\lambda(k)$  al-  
lows the search to evolve from an initial state where  
all target locations compete equally for a saccade to a  
final state where only a few most likely target loca-  
tions remain.

4. Move the eye to the location found by step (3). Al-  
though in our simulations we can get away with not  
actually implementing this step, as explained below.
5. Repeat steps (2), (3) and (4) above with  
 $k = \max - 1, \max - 2, \dots$  until either the target ob-  
ject has been foveated or the number of scales has  
been exhausted. In the former case, the decision pro-  
cess signals the termination of the search process. In  
the latter case, subsequent eye movements are made  
using saliency maps based on all the scales.

The model has only one parameter, the initial value of  
 $\lambda(1)$ . The function of  $\lambda(k)$  is to sharpen the peaks in the  
saliency map. The specific initial value of  $\lambda(1)$  is de-  
pendent on the values in the filter kernels. With each  
target computation,  $\lambda(k)$  was decreased by a factor of  
two, thereby allowing the search to evolve from an ini-  
tial coarse resolution state where many target correla-  
tions contribute to a saccade, to a final state where only  
a single most likely target location contributes. The  
values for  $\lambda(k)$  used were 4, 2 and 1 for  $k = 1, 2$  and 3  
respectively. The exact values are not crucial; the data  
can be fit qualitatively with values of  $\lambda(1)$  ranging from  
1 to 20. The same values of  $\lambda(k)$  are used for all scenes  
and target locations within a scene.

The modified targeting model was implemented on  
our pipeline image processor. Fig. 8 shows the saliency



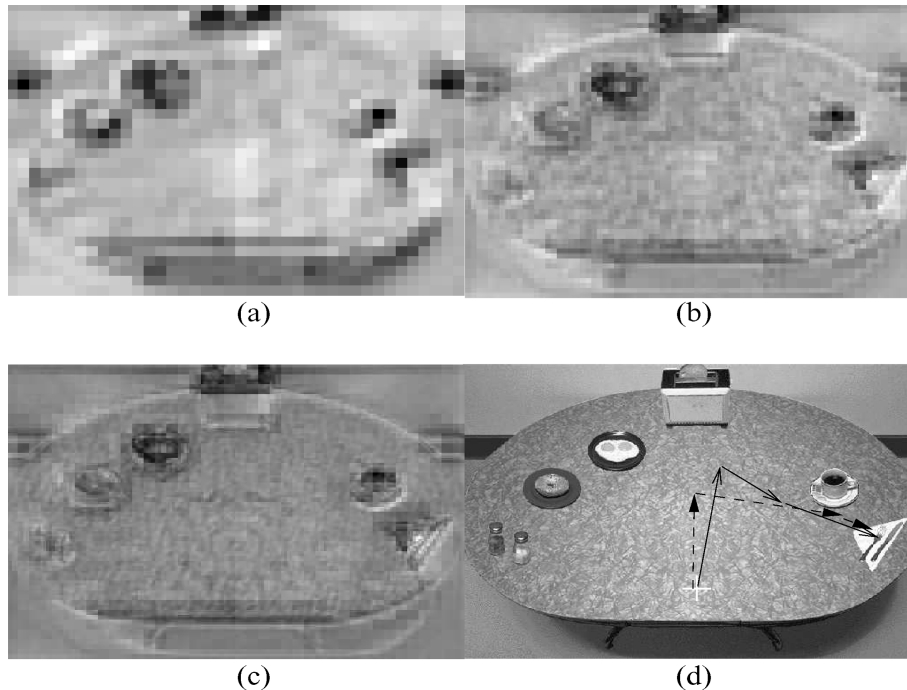


Fig. 8. Illustration of coarse-to-fine saccadic targeting. The saliency map  $S(x, y)$  after the inclusion of the largest (a), intermediate (b), and smallest scale (c) as given by filter response distances to the prototype (the fork and knife); the brightest points are the closest matches; (d) shows the predicted eye movements as determined by the weighted population averaging scheme. For comparison, saccades from a human subject are given by the dotted arrows.

533 maps for this image after each of three iterations, with  
534 the middle and highest frequencies included in (b) and  
535 (c) respectively. Part (d) of the figure shows the sequence  
536 of fixations generated by the model for this image, to-  
537 gether with those from a human subject. The target  
538 (composed of the fork and the knife) was the same in  
539 both cases. Thus the coarse-to-fine analysis, together  
540 with center-of-gravity effects, can produce the kind of  
541 fixation patterns that human subjects generate with this  
542 image.

543 In Fig. 8 the saliency map should of course be shifted  
544 with gaze. The reason we do not do this is simple ex-  
545 pediency. Since we assume the resolution is chosen at the  
546 outset of the search, this implies that it is not changed  
547 during the target selection, therefore the saliency map  
548 cannot take advantage of the resolution of the fovea  
549 during the targeting period. The reason for this may not  
550 be obvious: if the target is being decided upon by some  
551 kind of correlation, the correlation function for foveated  
552 targets and non-foveated targets must be adjusted in a  
553 way that depends on the eccentricity and target. Oth-  
554 erwise a false target near the fovea might appear better  
555 than an eccentric true target. This is avoided in the  
556 model by selecting a resolution based on the signal-to-  
557 noise properties of the display and using that resolution  
558 cutoff everywhere in the resultant search window. As a  
559 consequence the saliency map is, to a first approxima-

tion, just shifted by saccades. We do not shift it in our 560  
figures in order to more easily compare visually the 561  
temporal effect of sequentially applying the multiple- 562  
scale filters. 563

## 6. Model-data comparison

564

The model's performance was compared to human 565  
data taken with 480 search trials pooled over four sub- 566  
jects. Owing to the nature of the different distractors and 567  
targets, there is substantial intersubject variability for 568  
each configuration, nonetheless, on the average, the 569  
model is remarkably good at approximating the actual 570  
gaze changes that subjects make. To show this we did 571  
the following analyses. The first step was to separate the 572  
sequences that ended up on the target with those that 573  
went to neighboring targets. Over the 480 trials, many 574  
records showed eye movements to nearby targets. This 575  
data is consistent with observations of both Kowler and 576  
Findlay who showed, particularly in the case when eye 577  
movements are made immediately upon the onset of the 578  
display, that a percentage of the movements were to 579  
false targets. Interestingly, the model also makes eye 580  
movements to false targets, but generally not to the 581  
same ones made by the subjects. Thus to compare the 582  
two sets of data we did the following: 583

584 (1) We generated an *average observer's* path to each of  
585 the six locations by averaging the fixations over subjects  
586 and target images. The coordinates were weighted by the  
587 variance between subjects. This meant that if a subject's  
588 movements were dissimilar to the group, they counted  
589 less in the sum. In the small number of cases where there  
590 were more than three saccades, only the first three were  
591 counted, as by the third saccade the eyes were always  
592 very close to one of the targets.

593 (2) The model data was averaged over the different  
594 targets for each location. In addition, trials where the  
595 final saccade was closer to a false target were excluded  
596 from the data and scored as errors. This resulted in 27  
597 false targets in 120 model trials. In comparison, if we  
598 count human subject trials that had a standard deviation  
599 of the subjects' final gaze points of more than 75% of the  
600 intertarget separation difference as errors, then 29 of the  
601 records averaged over subjects are counted as false tar-  
602 gets.

603 After these steps the results are shown in Fig. 9. The  
604 box in each sub-figure represents a  $1^\circ$  region centered on  
605 each target location. As is evident there is very good  
606 agreement between the model and human data for each  
607 location. Furthermore the number of errors made by the  
608 model is in very close agreement with the number of  
609 errors made by human subjects. It would be perhaps  
610 desirable to have the model represent an average or  
611 prototypical subject, but we cannot do this as the filters  
612 used by the model are probably slightly different than  
613 those used by the subjects, as described subsequently.  
614 However, we can ask whether the model is representa-

615 tive of an individual subject, and there the evidence is  
616 very encouraging. The average standard deviation for  
617 the subjects, averaged overall fixations is  $1.5^\circ$  whereas  
618 the average difference between model and average sub-  
619 ject fixations is  $0.7^\circ$ . Thus the model behavior is well  
620 within the profile expected of an individual subject.

621 We also examined the saccades to false targets to see if  
622 there was any systematic bias in terms of location, target  
623 or scene type. One might well ask why there should be  
624 *any* false targets, since the decisions made by the sub-  
625 jects as to target presence are 100% accurate. We believe  
626 that the model provides an answer: (a) the decision  
627 process is separate from the targeting process and thus  
628 can still function when the ultimate target is eccentric,  
629 and (b) gaze can be mislocated since the template is  
630 defined on a neutral background and the background of  
631 the display bleeds into the larger filters, disturbing the  
632 correlation computation.

633 Table 1 shows this data for target location. The table  
634 shows the principal difference between the human and  
635 model data. The model had no difficulty with the crib  
636 scene, where targets were arrayed on a high contrast  
637 background, but the human subjects spread their errors  
638 around all three scenes uniformly. We interpret this to  
639 mean that the filter model is not identical to that used by  
640 the human subjects in that the filters are too sensitive to  
641 contrast and not sensitive enough to the fine structure in  
642 the targets. Nonetheless, given this caveat, the overall  
643 pattern of errors among locations is fairly uniform in  
644 both data sets.

645 Additional evidence for the correlation model comes  
646 from a control experiment that we performed, in which

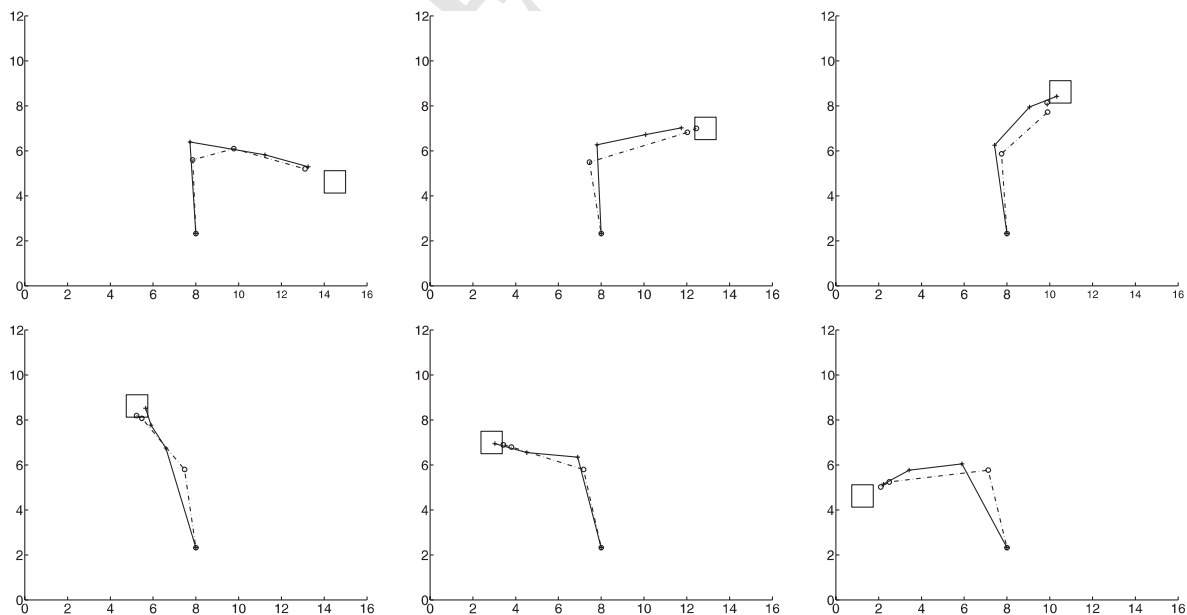


Fig. 9. Model vs human subjects results. The figure shows the performance of the subjects averaged over subjects and targets to each target location (see text). The scale is in degrees and the box shows a  $1^\circ$  region centered around each target. Circles and plus symbols mark the fixation points for human and model data respectively.

Table 1  
Number of false targets for the model and human subjects, broken down by target location and scene

Location	Crib	Dine	Work
<i>Model</i>			
1		3	3
2		3	1
3			3
4		2	
5		2	3
6		4	3
<i>Subjects</i>			
1	3	4	3
2	2	3	1
3		1	
4	3	2	2
5	1	1	
6		1	2

647 the contextual background (e.g. the workbench and  
648 other objects) was removed, and the search objects were  
649 presented on a uniform background. Table 2 shows the  
650 results in the form of initial endpoint error after the first  
651 saccade. A striking point of comparison is the difference  
652 in error for search scenes containing a single object in  
653 the case of a uniform color background (c) vs a non-  
654 uniform realistic background (a) and (b). In the former  
655 case, the error is reduced by a factor of two for color  
656 images and slightly more than that for the gray scale  
657 images. This result implies an interference due to the  
658 background in the targeting process, as assumed by the  
659 model. As one might expect, the effect of the back-  
660 ground is less as the number of target objects increases.  
661 This experiment is described in more detail in (Zelinsky  
662 et al., 1997). It is also of interest to compare the end-  
663 point error for color and gray scale images. A small  
664 difference is evident after the first saccade. After the  
665 second saccade, the endpoint error was a full 1° less in  
666 the case of color images, strongly suggesting that color  
667 information is being used in the targeting computation.  
668 Although the simulation results described in this section  
669 modeled human eye movement data from gray scale

Table 2  
The effect of background on saccade accuracy. Mean endpoint error (in degrees) across all four subjects after the first saccade as a function of three different display conditions: (a) color images with a realistic background, (b) gray scale images with a realistic background, and (c) color images with a uniform background

Condition	Set size		
	1	3	5
(a) Color	3.2	4.8	5.1
(b) Gray	3.8	5.0	5.2
(c) Uniform background	1.6	4.8	5.1

The errors are shown for set sizes of one, three, and five objects in the search scene. Note that a uniform background for one target causes initial saccade accuracy to increase by a factor of two, implying that the background and other targets are deviating the saccade trajectory.

images, the model can be readily extended for saccadic 670  
targeting based on color information. 671

## 7. Appearance-based search vs spatial memory search 672

In both the model and experiment there is no prior 673  
knowledge of the specific location of the target before 674  
the presentation of the search array. Thus the only in- 675  
formation available in both cases is *what* the target looks 676  
like, not *where* it is, and the search strategy is based 677  
primarily on the object's appearance. However, it seems 678  
intuitively likely that information about an object's lo- 679  
cation based on previous fixations in a continuously 680  
present scene, would contribute to the search process. 681  
Both physiological and psychophysical evidence reveal 682  
the ability to make saccades purely on the basis of in- 683  
formation about spatial location (Colby & Goldberg, 684  
1999). Precuing a location also reduces saccade latencies 685  
to that location. However, it is not clear what role 686  
spatial information plays when the stimulus is present 687  
on the retina and can be chosen on the basis of ap- 688  
pearance, as is ordinarily the case in natural viewing, 689  
where subjects have usually made multiple fixations in a 690  
scene. Evidence from natural tasks such as tapping 691  
(Epelboim, Steiman, Kowler, & Pizlo, 1997; Land, 692  
Mennie, & Rusted, 1999) suggest that spatial informa- 693  
tion does ordinarily play a role in the targeting process. 694  
Thus adding spatial information to the task should af- 695  
fect the targeting strategy. 696

To test whether spatial information in addition to 697  
appearance factors would change the search pattern, a 698  
modification of the visual search task described above 699  
was run, where subjects were allowed to briefly *preview* 700  
the search scene (without knowing the search target) in a 701  
separate interval just before the search target was pre- 702  
sented. Subjects were given a one second opportunity to 703  
preview the search scene prior to the presentation of the 704  
target. In this period, they were allowed to move their 705  
gaze freely, allowing them to fixate individual targets. 706  
The rest of the experiment remained the same as before 707  
(Zelinsky & Sheinberg, 1997). The subjects held fixation 708  
on a fixation cross, an icon of the target was then pre- 709  
sented at the fixation point, followed by the search 710  
scene. An analysis of the eye movement data revealed 711  
that single saccades were by far the most common, as 712  
summarized in Fig. 10. The histograms show the initial 713  
endpoint error after the first saccade for the original 714  
search paradigm and the same for the case where sub- 715  
jects had a one second preview of the scene containing 716  
the potential targets. For most but not all of the preview 717  
cases, the initial endpoint error is 1° or less, strongly 718  
suggesting that subjects use the spatial location of the 719  
targets as an integral part of the search process. In ad- 720  
dition, the reaction time for the decision was about 100 721  
ms faster when the preview was presented, suggesting 722

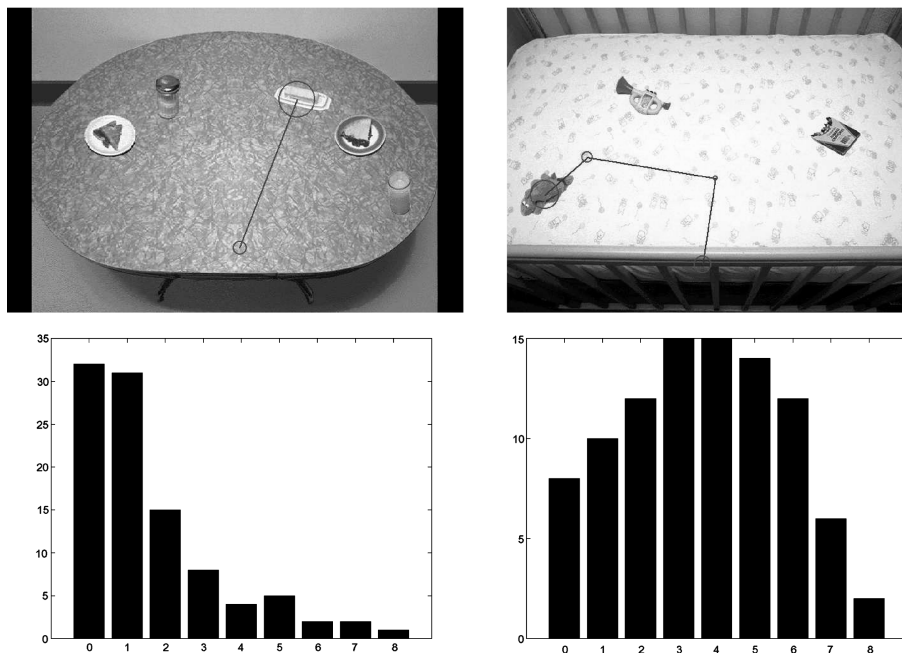


Fig. 10. Comparing preview vs no preview. The graph shows histograms of the endpoint error after the first saccade for the original search paradigm and the case where subjects had a one second preview of the potential targets. For most but not all of the preview cases the endpoint error is 1° or less, implying that subjects were able to remember and use the spatial location of the targets. Histograms: vertical axis = frequency of occurrences, horizontal axis = degrees.

723 that the location information facilitated the search  
724 process (Zelinsky & Sheinberg, 1997). This might occur  
725 if subjects were able to associate locations in the saliency  
726 map with the filter response vectors for objects, so that  
727 seeing one of these objects would now “prime” the  
728 corresponding location in the saliency map. This priming  
729 would in turn allow more accurate saccadic targeting  
730 in the cases where the target location happened to be  
731 inventoried during the preview period. It is important to  
732 remember that the subjects’ task was simply to respond  
733 with a key press whether the target was present or not.  
734 No instructions were made about eye movements except  
735 that the subject should fixate the cross before the stimulus  
736 presentation. Thus it is likely that the observers are  
737 integrating the spatial and appearance information as  
738 part of a natural search strategy that results in more  
739 direct saccades. A way to extend the model to do this is  
740 described in (Ballard et al., 1997).

## 741 8. Discussion

742 The current model shares some mechanisms used by  
743 Itti and Koch (2000). They also propose a specific  
744 computational implementation of stimulus saliency for  
745 general scenes. Itti and Koch also propose filtering the  
746 image at different spatial scales. However, their model  
747 differs in that separate saliency maps are computed for  
748 color, intensity, and orientation. These separate maps

are linearly combined following iterative lateral competition within each map. The saliency peak is then found using a winner-take-all network. Our model has a single saliency map by using oriented spatiochromatic filters, but the most important difference is that it uses a top-down search template to locate the saliency peak. Itti and Koch have no obvious way of searching for specific targets that are not contained in their bottom-up maps. Furthermore, our model works with unsegmented images, and thus avoids the difficult task of deciding what constitutes a “feature.” The other important difference is our evolution of the signal in time with the addition of information at higher spatial frequency which is needed to fit the human data. Itti and Koch also have no direct comparisons with human data.

The model used by Tsotsos (Tsotsos et al., 1995) is more similar to that described here in that it has a top-down target component. However, there is no attempt in the Tsotsos model to model the details of eye movements in a way that could capture the skipping saccades seen in human data.

The model shares some general similarities with the visual search model proposed (but not implemented) by Findlay and Walker (1999), as well as that of Hooge (1996). Their suggestion of a temporal evolution of the saliency map takes specific form here. We differ most from Findlay and Walker in the representation of temporal control. In our model there is no explicit temporal control of saccades other than the assumption that the

749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777

saliency map takes about 400 ms to evolve. We see this as a distinct biological advantage. By decoupling the dependence of the saliency map dynamics with the targeting system, they can be simpler and work independently.

Although not a computer model, the model used by Motter and Holsapple (2000) is very relevant to our work. Motter's studied monkeys search patterns in looking for small conjunction targets of color and shape and found that the data in different displays could be normalized by dividing by the density of search patterns of the correct color. He terms this an adjusted nearest neighbor distance (ANND). The reason this is relevant to our own model is that although not implemented, we conceptualize the search window as being adjusted based on signal-to-noise characteristics. The ANND concept can be seen as making a similar suggestion as dense target arrays can reduce signal-to-noise as shown by Palmer et al. (2000).

The model also shares some similarities with that of Wolfe, Cave, and Franzel (1989) and might be seen as an extension that fixes important problems with that model. In the Wolfe and Cave model, top-down priming of features in the saliency map computations can direct the search. Important differences arise in how these computations are carried out. To implement these calculations, their model requires that the features be segmented from the background, an unrealistic requirement in general. In contrast, our general correlation-based targeting method can handle arbitrary targets. More importantly, by separating the eye movements from the decision process, as is done in our model, means that gaze does not have to search every item in a multiple-item search task, but can use area-based calculations. The skipping data provides evidence that this can happen as the eyes move to non-target locations en route to making a decision. Motter's ANND data and Zelinsky's data provide further evidence for area-based vs item-by-item search.

Explaining the observed skipping saccades is done using a coarse-to-fine matching mechanism. The main benefit of a coarse-to-fine strategy is that it allows continuous execution of the decision and oculomotor processes, thereby increasing the probability of an early match. Coarse-to-fine strategies have also enjoyed recent popularity in computer vision with the advent of image pyramids for tasks such as motion detection (Burt, 1988). One key question that remains is the source of sequential application of the filters in the human visual system. This will usually result from the variation in resolution of the retina. Since resolution falls off with distance from the fovea, the fine spatial scales could be ineffective during early stages of search simply because the fixation point is distant from the target. However, our model suggests a different explanation. First, the three filters used in the model predictions were centered

about 1, 3, and 6 cycles per degree. Even the highest of these should be visible at an eccentricity of  $7^\circ$  (Anderson et al., 1991). To test if the targets were identifiable at this eccentricity, in a control experiment observers were required to identify the targets while maintaining fixation. They were able to do this with negligible errors but used much longer reaction times (Zelinsky & Sheinberg, 1997). In addition, in the experiment where subjects were given a preview, many saccades went directly to the target, suggesting that resolution did not preclude direct targeting. Since the model fits the data well, it suggests that the additional effects on targeting from higher acuity measurements might be small.

An additional explanation for the sequential application of the filters is that the cortical machinery is setup to match the larger scales first, as target information is propagated via cortico-cortical feedback from higher to lower areas in the visual cortical hierarchy. If this were the case, the observed data would result from the fact that the oculomotor system is ready to move before all the scales can be matched, and thus the eyes move to the current best target position. This interpretation of the data is appealing for two reasons. First, it reflects a long history of observations on the priority of large scale channels in vision (Breitmeyer, 1975; Navon, 1977; Parker & Dutch, 1987). A particularly relevant experiment is that of Schyns and Oliva (1994). This shows that in a recognition task with 30 ms exposures, subjects are sensitive to the low frequencies in the image whereas with 150 ms exposures, subjects respond to the high frequency content. Second, in a search experiment similar to ours done by Findlay (1997), when subjects held their gaze before starting the search, the pattern of saccades was more direct, suggesting that the target location had been refined during the wait. In another experiment using pairs of targets, Findlay (1997) found evidence that the saccade target signal is initially coarsely localized, and becomes more refined with increasing duration. Thus it is not clear whether the coarse-to-fine analysis is instantiated in the hardware or whether it is a de facto consequence of peripheral resolution fall off. Even if peripheral information is not limiting in a particular instance, coarse-to-fine analysis may develop as a naturally efficient strategy, since foveation will invariably lead to additional high frequency information for the current perceptual decision.

An alternative explanation for the initial saccade towards the center of the display is that it is a preplanned saccade to facilitate the search by centering fixation within the search array. The brief latencies before the first saccade support the idea of some kind of preprogramming. However, it is not likely to be entirely strategic (as opposed to a center-of-gravity saccade) because the initial fixation is biased toward the target.

One might suspect that the findings were a product of the experimental setup, which had subjects's heads fixed

834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889

890 in a bite-bar. To check this we repeated the tests using a  
891 stereoview head mounted display which contained an  
892 eye tracker. We did not analyze the results quantita-  
893 tively, but skipping movements were ubiquitous in the  
894 data.

895 Normally, a saccade is followed by a 200–300 ms  
896 fixation period before the next saccade is generated.  
897 Under certain circumstances, *express saccades* are also  
898 observed (Fischer & Boch, 1983; Fischer & Ramsperger,  
899 1984; Fischer & Weber, 1993). The fixation periods for  
900 express saccades are much shorter, in the range 70–100  
901 ms. An analysis of the visual search results (Zelinsky et  
902 al., 1997) revealed that the fixation periods of some of  
903 the center-of-gravity “skipping” eye movements are  
904 much smaller than normal (around 80–130 ms), small  
905 enough to qualify them as express saccades. There is a  
906 very simple explanation of these short latencies in the  
907 context of the proposed model. In a normal fixation,  
908 information from that fixation is presumably used in the  
909 computation of the next target. This necessitates some  
910 setup time for the information to be part of the targeting  
911 computation. However, in some cases, the next target  
912 may not require information from the current fixation.  
913 In such cases, the fixation times can be made much  
914 shorter. Such a situation may occur in the case of the  
915 “skipping” eye movements, as the targeting is based on  
916 a correlation process which is being done sequentially  
917 across scales. Of course, the partial correlation results  
918 contained in the saliency map have to be “shifted” due  
919 to the intermediate eye movements, before being inte-  
920 grated, but the eye movement itself contains the infor-  
921 mation necessary to perform this shifting. The crucial  
922 point is that express saccades may simply reflect a simple  
923 relationship between the ongoing computation of the  
924 saliency map and the motor command that executes eye  
925 movements. When the saliency map computations can  
926 be speeded up, the rate of saccades can be made corre-  
927 spondingly faster.

928 There exists a vast literature on the role of attention in  
929 visual cognition (Duncan & Humphreys, 1992; Krose &  
930 Julesz, 1989; Posner & Petersen, 1990; Saarinen & Ju-  
931 lesz, 1991; Treisman, 1988; Treisman & Gelade, 1980).  
932 Attention has been characterized as covert search based  
933 on the metaphor of an attentional spotlight. Some of the  
934 search results have suggested that targets can be exam-  
935 ined at the rate of about 25 ms per item, with the at-  
936 tentional spotlight moving from one location to the next  
937 at a speed of about one attentional shift every 30–50 ms  
938 (Krose & Julesz, 1989; Saarinen & Julesz, 1991). Models  
939 of attention (for example, Niebur & Koch, 1996) have in  
940 fact literally modeled this shift of the “focus of atten-  
941 tion”. The technical advantage of such a strategy is that,  
942 since gaze is fixed, retinal coordinates can be used for  
943 keeping track of examined locations. However, since  
944 signal transmission through visual cortex is on the order  
945 of 80–100 ms, performing covert search with an atten-

tional spotlight while simultaneously obeying this 946  
stringent time constraint seems a difficult endeavor. An 947  
alternate explanation provided by the present model is 948  
that covert search occurs whenever the decision process 949  
finishes before an eye movement is made. This would 950  
occur, for example, in the cases where the presence of 951  
the target in a peripheral location can be judged directly 952  
from the correlation peaks in the saliency map using a 953  
signal-to-noise criterion. Under such circumstances, the 954  
eye movement becomes superfluous and a decision as to 955  
the presence or absence of the target can be made im- 956  
mediately without the need for an overt saccade. Such 957  
an interpretation is especially attractive since it allows a 958  
single targeting mechanism to parsimoniously account 959  
for both covert and overt search. It is also consistent 960  
with a body of evidence suggesting that the “atten- 961  
tional” (decision-making) and saccadic systems are 962  
regulated by different but closely related oculomotor 963  
control systems (Shepherd, Findlay, & Hockey, 1986; 964  
Groner, 1988; Corbetta, 1999; Findlay, 1997; Motter & 965  
Belky, 1998; Rizzolatti, 1996). The model has the addi- 966  
tional advantage of being simpler than models that use 967  
additional machinery to couple the decision and tar- 968  
geting systems (e.g. Findlay, 1997). 969

## 9. Conclusion 970

A large number of computational models pertaining 971  
to human visual search and attention have previously 972  
been proposed (Chapman, 1991; Niebur & Koch, 1996; 973  
Olshausen, Van Essen, & Anderson, 1993; Tsioutsias & 974  
Mjolsness, 1996; Tsotsos et al., 1995; Wolfe, 1994). 975  
Many of these rely on predominantly bottom-up atten- 976  
tional processes based on various forms of feature maps 977  
that are used to facilitate search. Some of these models 978  
were motivated primarily by the need to explain classical 979  
reaction time results rather than the pattern of eye 980  
movements observed during visual tasks. Other models 981  
have explored the use of bottom-up saliency maps and 982  
have used eye movement scan-paths as sensorimotor 983  
memories for recognition (Didday & Arbib, 1975; Gie- 984  
fing, Janßen, & Mallot, 1991; Rimey & Brown, 1991; 985  
Rybak, Gusakova, Golovan, Podladchikova, & Shev- 986  
tsova, submitted for publication; Yamada & Cottrell, 987  
1995). This paper proposes a new model of the gaze 988  
targeting process in natural tasks based on observations 989  
of (Geisler & Chou, 1995; Motter & Holsapple, 2000; 990  
Palmer et al., 2000) that uses both bottom-up scene 991  
representations as well as top-down target descriptions 992  
for gaze control. 993

The model has four principal features: 994

(1) Instead of “features” that are preselected inde- 995  
pendently of a task, the model uses iconic templates that 996  
are task-dependent. As they are expressed in terms of 997

image filter responses, that are both more general and simpler to use than features. Eye movement models that are based on a fixed library of features cannot explain how arbitrary targets are computed.

(2) The model separates the process of changing gaze from that of deciding on properties of a target. This has the virtue of allowing the timing relationships between these two processes to be a natural consequence of the properties of the scene. This greatly simplifies the control problem of coordinating eye movements and decisions.

(3) The model specifies that the correlation used to select targets proceeds in a coarse-to-fine manner that takes time. If the target is novel and its location must be determined solely on appearance, this time is longer than that needed to generate an eye movement, and consequently effects the gaze trajectory in a predictable way. This result provides a concrete model of a myriad of experimentally observed "center-of-gravity" observations. Since our center of gravity is correlation-based, it is readily tested experimentally.

(4) The most controversial aspect of the model is its use of area-based search. The assumption is that the resolution used to search for the target can be chosen at the beginning of the search based on the signal-to-noise properties of the search area. The motivation for being able to do this is to search large areas at comparable resolution. The assumption that humans would not make continuous use of all the available resolution in the retinotopic array is counterintuitive. We have argued that it has precedents in search models, and our experiments show (1) that the model fits the data well and (2) foveal resolution is not necessary for target location. However we cannot rule out the use of all the available resolution by human subjects, so that this question needs to be settled by further experiments.

The model is constructive, has a specific computational prescription for target computation, and fits experimental observations. Its most controversial claim is that, for the experimental conditions tested, it can use resolutions much lower than that ultimately available from the scene to guide gaze changes. As a consequence, the effect of additional foveal resolution has minimal effects on the gaze trajectory. We anticipate that situations could be constructed for which foveal effects would be seen, but those effects may prove a refinement on the model presented here.

The main goal of the model was to capture the exogenous effects of the visual stimulus. There has been no attempt to model endogenous target specifications e.g. anti-saccades. However these effects have been modeled by Kopecz and Schoner (1995) and Trappenberg, Dorris, Munoz, and Klein (2001) in a way that is compatible with our model.

## Acknowledgements

This work was supported by NSF research grant no. CDA-8822724, NIH/PHS research grants no. 1-R24-RRO6853, EY-05729 and 1-P41-RR09283, and a grant from the Human Science Frontiers Program. The paper greatly benefitted from the reviewers' comments.

## References

- Anderson, S. J., Mullen, K. T., & Hess, R. F. (1991). Human peripheral resolution for chromatic and achromatic stimuli-limits imposed by optical and retinal factors. *Vision Research*, 442, 47-64.
- Ballard, D. H., Hayhoe, M., & Pelz, J. (1995). Memory representations in natural tasks. *Journal of Cognitive Neuroscience*, 7(1), 66-80.
- Ballard, D. H., Hayhoe, M. M., Pook, P. K., & Rao, R. P. N. (1997). Deictic codes for the embodiment of cognition. *Behavioral and Brain Sciences*, 20.
- Barrow, H. G. (1987). Learning receptive fields. In *Proceedings of the IEEE International Conference on Neural Networks* (pp. 115-121).
- Bell, A. J., & Sejnowski, T. J. (1997). The 'independent components' of natural scenes are edge filters. *Vision Research*.
- Bichot, M., & Schall, J. (1999). Saccadic target selection in the macaque during feature and conjunction visual search. *Visual Neuroscience*, 16, 81-89.
- Breitmeyer, B. G. (1975). Simple reaction time as a measure of the temporal response properties of transient and sustained channels. *Vision Research*, 15, 1411-1412.
- Bulthoff, H., & Edelman, S. (1992). Psychophysical support for a two-dimensional view interpolation theory of object recognition. *Proceedings of the National Academy of Science USA*, 89, 60-64.
- Burt, P. J. (1988). Attention mechanisms for vision in a dynamic world. In *ICPR* (pp. 977-987).
- Chapman, D. (1991). *Vision, instruction, and action*. Cambridge, MA: MIT Press.
- Colby, C., & Goldberg, M. (1999). *Annual Review of Neuroscience*, 22, 319-349.
- Coren, S., & Hoenig, P. (1972). Effect of non-target stimuli upon length of voluntary saccades. *Perceptual and Motor Skills*, 34, 499-508.
- Derrico, J. B., & Buchsbaum, G. (1991). A computational model of spa-tiochromatic image coding in early vision. *Journal of Visual Communication and Image Representation*, 2(1), 31-38.
- Didday, R. L., & Arbib, M. A. (1975). Eye movements and visual perception: A two visual system model. *International Journal of Man-Machine Studies*, 7, 547-569.
- Duncan, J., & Humphreys, G. W. (1992). Beyond the search surface: visual search and attentional engagement. *Journal of Experimental Psychology: Human Perception and Performance*, 18, 578-588.
- Eckstein, Miguel P. (1998). The lower visual search efficiency for conjunctions is due to noise not serial attentional processing. *Psychological Science*, 9, 111-118.
- Epelboim, J., Steiman, R. M., Kowler, E., & Pizlo, Z. (1997). Gaze-shift dynamics in two kinds of sequential looking tasks. *Vision Research*, 37, 2597.
- Findlay, J. (1982). Global visual processing for saccadic eye movements. *Vision Research*, 22, 1033-1045.
- Findlay, J. (1987). Visual computation and saccadic eye movements: A theoretical perspective. *Spatial Vision*, 2, 175-189.
- Findlay, J. (1997). Saccade target selection during visual search. *Vision Research*, 37, 617-631.
- Fischer, B., & Boch, R. (1983). Saccadic eye movements after extremely short reaction times in the monkey. *Brain Research*, 260, 21-26.

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

- 1113 Fischer, B., & Ramsperger, E. (1984). Human express-saccades:  
1114 extremely short reaction times of goal directed eye movements.  
1115 *Experimental Brain Research*, 57, 191-195.
- 1116 Fischer, B., & Weber, H. (1993). Express saccades and visual attention.  
1117 *Behavioral and Brain Sciences*, 16, 553-610.
- 1118 Freeman, W. T., & Adelson, E. H. (1991). The design and use of  
1119 steerable filters. *IEEE Transactions on Pattern Analysis and*  
1120 *Machine Intelligence*, 13(9), 891-906.
- 1121 Geisler, W. S., & Chou, K.-L. (1995). Separation of low-level and high-  
1122 level factors in complex tasks: visual search. *Psychological Review*,  
1123 102(2), 356-378.
- 1124 Gieffing, G.-J., Janßen, H., & Mallot, H. (1991). A saccadic camera  
1125 movement system for object recognition. In T. Kohonen, K.  
1126 Makisara, O. Simula, & J. Kangas (Eds.), *Artificial Neural*  
1127 *Networks* (vol. 1) (pp. 63-68). Amsterdam: Elsevier.
- 1128 Gilchrist, I., & Heywood, C. (1999). Saccade selection in visual search:  
1129 evidence for spatial frequency specific between item interactions.  
1130 *Vision Research*, 39, 1373-1393.
- 1131 Graham, N. (1989). *Visual pattern analyzers*. New York: Oxford  
1132 University Press.
- 1133 Groner, R. (1988). Eye movements, attention and visual information  
1134 processing: some experimental results and methodological consid-  
1135 erations. In G. Luer, U. Lass, & J. Shallo-Hoffman (Eds.), *Eye*  
1136 *Movement Research: Physiological and Psychological Aspects* (pp.  
1137 295-319). Göttingen, Germany: Hogrefe.
- 1138 Hancock, P. J. B., Baddeley, R. J., & Smith, L. S. (1992). The principal  
1139 components of natural images. *Network*, 3, 61-70.
- 1140 He, P., & Kowler, E. (1989). The role of location probability in the  
1141 programming of saccades: Implications for "center-of-gravity"  
1142 tendencies. *Vision Research*, 29, 1165-1181.
- 1143 Hooge, I., & Erkelens, C. (1998). Adjustment of fixation duration  
1144 during visual search. *Vision Research*, 38, 1295-1302.
- 1145 Hooge, I. T. C. (1996). *Control of eye movements in visual search*. Ph.D.  
1146 Thesis. Netherlands: University of Utrecht.
- 1147 Itti, L., & Koch, C. (2000). A saliency-based search mechanism for  
1148 overt and covert shifts of visual attention. *Vision Research*, 40, 11-  
1149 46.
- 1150 Just, M. A., & Carpenter, P. A. (1976). Eye fixations and cognitive  
1151 processes. *Cognitive Psychology*, 8, 441-480.
- 1152 Kanerva, P. (1988). *Sparse distributed memory*. Cambridge, MA:  
1153 Bradford Books.
- 1154 Koch, C., & Ullman, S. (1985). Shifts in selective visual attention:  
1155 Toward the underlying neural circuitry. *Human Neurobiology*, 4(4),  
1156 219-227.
- 1157 Kopecz, K., & Schoner, G. (1995). Saccadic motor planning by  
1158 integrating visual information and pre-information on neural  
1159 dynamic fields. *Biological Cybernetics*, 73, 49-60.
- 1160 Kowler, E., & Anton, S. (1987). Reading twisted text: implications for  
1161 the role of saccades. *Vision Research*, 27, 45-60.
- 1162 Krose, B. J. A., & Julesz, B. (1989). The control and speed of shifts of  
1163 attention. *Vision Research*, 29(11), 1607-1619.
- 1164 Land, M. F., & Furneaux, S. (1997). The knowledge base of the  
1165 oculomotor system. In *Proceedings of the Royal Society Conference*  
1166 *on Knowledge-Based Vision*, February.
- 1167 Land, M., Mennie, M., & Rusted, J. (1999). An active role of  
1168 vision and eye movements in the control of activities of daily living.  
1169 *Perception*, 28, 1311-1328.
- 1170 McPeck, R., & Keller, E. (2001). Short term priming, concurrent  
1171 processing and saccade curvature during a target selection task in  
1172 the monkey. *Vision Research*, 41, 785-800.
- 1173 Motter, B., & Belky, E. (1998). The guidance of eye movements during  
1174 active visual search. *Vision Research*, 38, 1805-1815.
- 1175 Motter, B. C., & Holsapple, J. W. (2000). Cortical image density  
1176 determines the probability of target discovery during active search.  
1177 *Vision Research*, 40, 1311.
- 1178 Navon, D. (1977). Forest before trees: the precedence of global  
1179 features in visual perception. *Cognitive Psychology*, 9, 353-383.
- Niebur, E., & Koch, C. (1996). Control of selective visual attention:  
1180 modeling the "where" pathway. In D. Touretzky, M. Mozer, & M.  
1181 Hasselmo (Eds.), *Advances in Neural Information Processing*  
1182 *Systems* (vol. 8) (pp. 802-808). Cambridge, MA: MIT Press.
- 1183 Nowlan, S. J. (1990). Maximum likelihood competitive learning. In D.  
1184 S. Touretzky (Ed.), *Advances in Neural Information Processing*  
1185 *Systems* (vol. 2) (pp. 574-582). Morgan Kaufmann: San Mateo,  
1186 CA.
- 1187 Olshausen, B. A., & Field, D. J. (1996). Emergence of simple-cell  
1188 receptive field properties by learning a sparse code for natural  
1189 images. *Nature*, 381, 607-609.
- 1190 Olshausen, B. A., Van Essen, D. C., & Anderson, C. H. (1993). A  
1191 neurobiological model of visual attention and invariant pattern  
1192 recognition based on dynamic routing of information. *Journal of*  
1193 *Neuroscience*, 13, 4700-4719.
- 1194 Oram, M. W., & Perrett, D. I. (1992). Time course of neural responses  
1195 discriminating different views of the face and head. *Journal of*  
1196 *Neurophysiology*, 68(1), 70-84.
- 1197 O'Regan, J. K. (1990). Eye movements and reading. In E. Kowler  
1198 (Ed.), *Eye Movements and Their Role in Visual and Cognitive*  
1199 *Processes* (pp. 455-477). New York: Elsevier.
- 1200 Palmer, J., Vergese, P., & Pavel, M. (2000). The psychophysics of  
1201 visual search. *Vision Research*, 40, 1227.
- 1202 Parker, D. M., & Dutch, S. (1987). Perceptual latency and spatial  
1203 frequency. *Vision Research*, 27, 1279-1283.
- 1204 Poetzsch, M., Krueger, N., & Von der Malsburg, C. (1996). Improving  
1205 object recognition by transforming gabor filter responses. *Network*,  
1206 11, 341.
- 1207 Posner, M. I., & Petersen, S. E. (1990). The attention system of the  
1208 human brain. *Annual Review of Neuroscience*, 13, 25-42.
- 1209 Rao, R. P. N., & Ballard, D. H. (1995a). An active vision architecture  
1210 based on iconic representations. *Artificial Intelligence (Special*  
1211 *Issue on Vision)*, 78, 461-505.
- 1212 Rao, R. P. N., & Ballard, D. H. (1995b). Learning saccadic eye  
1213 movements using multiscale spatial filters. In G. Tesauro, D. S.  
1214 Touretzky, & T. K. Leen (Eds.), *Advances in Neural Information*  
1215 *Processing Systems* (vol. 7) (pp. 893-900). Cambridge, MA: MIT  
1216 Press.
- 1217 Rao, R. P. N., & Ballard, D. H. (1996). Dynamic model of visual  
1218 recognition predicts neural response properties in the visual cortex.  
1219 *Neural Computation*, 9, 721-763.
- 1220 Rao, R. P. N., & Ballard, D. H. (1997). Dynamic model of visual  
1221 recognition predicts neural response properties in the visual cortex.  
1222 *Neural Computation*, 9(4), 721-763.
- 1223 Rimey, R. D., & Brown, C. M. (1991). Controlling eye movements  
1224 with hidden Markov models. *International Journal of Computer*  
1225 *Vision*, 7(1), 47-65.
- 1226 Rybak, L. A., Gusakova, V. I., Golovan, A. V., Podladchikova, L. N.,  
1227 & Shevtsova, N. A. (submitted for publication). A model of attention-  
1228 guided visual perception and recognition. *Vision Research*.
- 1229 Saarinen, J., & Julesz, B. (1991). The speed of attentional shifts in the  
1230 visual field. *Proceedings of the National Academy of Science, USA*,  
1231 88, 1812-1814.
- 1232 Schyns, P. G., & Oliva, A. (1994). From blobs to edges: evidence for  
1233 time and spatial scale dependent scene recognition. *Psychological*  
1234 *Science*, 5, 195-200.
- 1235 Shepherd, M., Findlay, J., & Hockey, R. (1986). The relationship  
1236 between eye movements and spatial attention. *Quarterly Journal of*  
1237 *Experimental Psychology*, 38A, 475-491.
- 1238 Sperling, G., & Doshier, B. A. (1986). The attention operating  
1239 characteristic: some examples from visual search. *Science*, 202,  
1240 315-318.
- 1241 Tanenhaus, M. K., Spivey-Knowlton, M. J., Eberhard, K. M., &  
1242 Sedivy, J. E. (1995). Integration of visual and linguistic information  
1243 in spoken language comprehension. *Science*, 268, 632-634.
- 1244 Thorpe, S. J., & Imbert, M. (1989). Biological constraints on connec-  
1245 tionist modelling. In R. Pfeifer, Z. Schreier, F. Fogelman-Soulie, &  
1246



- 1247 L. Steels (Eds.), *Connectionism in Perspective* (pp. 63–92). Amsterdam: Elsevier. 1269
- 1248
- 1249 Trappenberg, T. P., Dorris, M. C., Munoz, D. P., & Klein, R. M. (2001). A model of saccade initiation based on the competitive 1270
- 1250 integration of exogenous and endogenous signals in the superior 1271
- 1251 colliculus. *Journal of Cognitive Neuroscience*, 13, 256–271. 1272
- 1252
- 1253 Treisman, A. (1988). Features and objects the fourteenth Bartlett 1273
- 1254 memorial lecture. *The Quarterly Journal of Experimental Psychology*, 40(2), 201–237. 1274
- 1255
- 1256 Treisman, A., & Gelade, G. (1980). A feature-integration theory of 1275
- 1257 attention. *Cognitive Psychology*, 12, 97–136. 1276
- 1258
- 1259 Tsioutsias, D. I., & Mjolsness, E. (1996). A multiscale attentional 1277
- 1260 framework for relaxation neural networks. In D. Touretzky, M. 1278
- 1261 Mozer, & M. Hasselmo (Eds.), *Advances in Neural Information 1279*
- 1262 Processing Systems (Vol. 8) (pp. 633–639). Cambridge, MA: MIT 1280
- 1263 Press. 1281
- 1264
- 1265 Tsotsos, J. K., Culhane, S. M., Wai, W. Y. K., Lai, Y., Davis, N., & 1282
- 1266 Nufflo, F. (1995). Modeling visual attention via selective tuning. 1283
- 1267 *Artificial Intelligence (Special Issue on Vision)*, 78, 507–545. 1284
- 1268
- 1269 Viviani, P. (1990). Eye movements in visual search: cognitive, 1285
- 1270 perceptual and motor control aspects. In *Eye Movements and 1286*
- 1271 Their Role in Visual and Cognitive Processes, New York. 1287
- 1272
- 1273 Weber, J., & Malik, J. (1995). Robust computation of optic flow in a 1288
- 1274 multiscale differential framework. *International Journal of Com- 1289*
- 1275 puter Vision, 14, 67–81. 1290
- 1276
- 1277 Wilson, H. R., & Wilkinson, F. (1997). Evolving concepts of spatial 1291
- 1278 channels in vision: from independence to nonlinear interactions. 1292
- 1279 *Perception*, 26, 939–960. 1293
- 1280
- 1281 Wolfe, J., Cave, K., & Franzel, S. (1989). Guided search: an alternative 1294
- 1282 to the feature integration model for visual search. *Journal of 1295*
- 1283 Experimental Psychology: Human Perception and Performance, 15, 1296
- 1284 419–433. 1297
- 1285
- 1286 Wolfe, J. (1994). Visual search in continuous naturalistic stimuli. 1298
- 1287 *Vision Research*, 34, 1187–1195. 1299
- 1288
- 1289 Yamada, K., & Cottrell, G. W. (1995). A model of scan paths applied 1300
- 1290 to face recognition. In *Proceedings of the 17th Annual Conference of 1301*
- 1291 the Cognitive Science Society (pp. 55–60). 1302
- 1292
- 1293 Yarbus, A. (1967). *Eye movements and vision*. New York: Plenum 1303
- 1294 Press. 1304
- 1295
- 1296 Zelinsky, G., & Sheinberg, D. (1997). Eye movements during parallel- 1305
- 1297 serial visual search. *Journal of Experimental Psychology: Human 1306*
- 1298 Perception and Performance, 23, 244–262. 1307
- 1299
- 1300 Zelinsky, G. J., Rao, R. P. N., Hayhoe, M. M., & Ballard, D. H. 1308
- 1301 (1997). Eye movements reveal the spatio-temporal dynamics of 1309
- 1302 visual search. *Psychological Science*. 1310