

# Online Research @ Cardiff

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository: <https://orca.cardiff.ac.uk/id/eprint/153272/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Pedziwiatr, Marek A., von dem Hagen, Elisabeth ORCID:  
<https://orcid.org/0000-0003-1056-8196> and Teufel, Christoph ORCID:  
<https://orcid.org/0000-0003-3915-9716> 2022. Knowledge-driven perceptual organization reshapes information sampling via eye movements. Journal of Experimental Psychology: Human Perception and Performance file

Publishers page:

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies.

See

<http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24

**Knowledge-driven perceptual organization reshapes information sampling via eye movements**

Marek A. Pedziwiatr <sup>1,2</sup>, Elisabeth von dem Hagen <sup>3</sup>, and Christoph Teufel <sup>1</sup>

<sup>1</sup> Biological and Computational Vision Lab, Cardiff University Brain Research Imaging Centre (CUBRIC), School of Psychology, Cardiff University

<sup>2</sup> Department of Biological and Experimental Psychology, Queen Mary University of London

<sup>3</sup> Cardiff University Brain Research Imaging Centre (CUBRIC), School of Psychology, Cardiff University

Word count: 13,728

**Author Note**

Marek A. Pedziwiatr  <https://orcid.org/0000-0002-3959-8666>

Elisabeth von dem Hagen  <https://orcid.org/0000-0003-1056-8196>

Christoph Teufel  <https://orcid.org/0000-0003-3915-9716>

We have no conflicts of interest to disclose.

Data from this study is openly available under the following link: [link to be provided upon publication]

Correspondence concerning this article should be addressed to Marek A. Pedziwiatr and Christoph Teufel, Cardiff University Brain Research Imaging Centre (CUBRIC), Cardiff University, Maindy Road, Cardiff, Wales, UK, CF24 4HQ. Emails: [m.pedziwiatr@qmul.ac.uk](mailto:m.pedziwiatr@qmul.ac.uk) and [teufelc@cardiff.ac.uk](mailto:teufelc@cardiff.ac.uk)

25

**Abstract**

26 Humans constantly move their eyes to explore the environment. However, how image-  
27 computable features and object representations contribute to eye-movement control is an ongoing  
28 debate. Recent developments in object perception indicate a complex relationship between  
29 features and object representations, where image-independent object-knowledge generates  
30 objecthood by reconfiguring how feature space is carved up. Here, we adopt this emerging  
31 perspective, asking whether object-oriented eye-movements result from gaze being guided by  
32 image-computable features, or by the fact that these features are bound into an object  
33 representation. We recorded eye movements in response to stimuli that initially appear as  
34 meaningless patches but are experienced as coherent objects once relevant object-knowledge  
35 has been acquired. We demonstrate that fixations on identical images are more object-centred,  
36 less dispersed, and more consistent across observers once these images are organised into  
37 objects. Gaze guidance also showed a shift from exploratory information-sampling to exploitation  
38 of object-related image-areas. These effects were evident from the first fixations onwards.  
39 Importantly, eye-movements were not fully determined by knowledge-dependent object  
40 representations but were best explained by the integration of these representations with image-  
41 computable features. Overall, the results show how information sampling via eye-movements is  
42 guided by a dynamic interaction between image-computable features and knowledge-driven  
43 perceptual organization.

44

45 *Keywords:* eye movements; perceptual organization; prior knowledge; object perception; natural  
46 scenes

47 *Public Significance Statement:* To explore and make sense of the world around us, we have to  
48 move our eyes. This study shows how our brain combines simple image-features such as edges  
49 and contrast with knowledge about objects to guide our eyes through a visual scene.

50

**51 Knowledge-driven perceptual organization reshapes information****52 sampling via eye movements**

53 Human visual experience carves up the world into objects (Feldman, 2003; Wagemans  
54 et al., 2012), distinct entities that are critical in structuring our interaction with the environment.  
55 When searching for a specific item in a scene or when exploring the world with no purpose  
56 other than to obtain information, humans tend to look at the centre of objects (e.g., Nuthmann  
57 & Henderson, 2010; Pajak & Nuthmann, 2013; Stoll, Thrun, Nuthmann, & Einhäuser, 2015).  
58 While these object-oriented effects of information sampling are well established, the current  
59 literature provides little consensus about which specific aspects of objects influence  
60 programming of eye movements (Borji & Tanner, 2016; Federico & Brandimonte, 2019; Hayes  
61 & Henderson, 2021; Henderson, Malcolm, & Schandl, 2009; Kilpeläinen & Georgeson, 2018;  
62 Nuthmann, Schütz, & Einhäuser, 2020; Van der Linden, Mathôt, & Vitu, 2015). This issue is  
63 complicated by the fact that it is often not clear exactly what constitutes an ‘object’ or how  
64 objects relate to image-computable features: except for special cases such as hallucinations  
65 (Horga & Abi-Dargham, 2019; Powers, Mathys, & Corlett, 2017; Teufel et al., 2015), features  
66 are necessary for visual object representations to arise but they are often not sufficient.  
67 Indeed, a growing number of studies using human psychophysics (Christensen, Bex, & Fiser,  
68 2015; Lengyel, Nagy, & Fiser, 2021; Lengyel et al., 2019; Neri, 2017; Ongchoco & Scholl,  
69 2019; Teufel, Dakin, & Fletcher, 2018) neuroimaging (Flounders, González-García, Hardstone,  
70 & He, 2019; Hsieh, Vul, & Kanwisher, 2010) , and animal electrophysiology (Gilbert & Li, 2013;  
71 Liang et al., 2017; Self et al., 2019; Self, van Kerkoerle, Supèr, & Roelfsema, 2013; Walsh,  
72 McGovern, Clark, & O’Connell, 2020) suggest that in order for object representations to  
73 emerge, prior object-knowledge has to interact with sensory processing. By contrast to  
74 conventional models of object recognition (DiCarlo, Zoccolan, & Rust, 2012; Kourtzi & Connor,  
75 2011; Kriegeskorte, 2015; Marr & Nishihara, 1978), these studies demonstrate that prior  
76 object-knowledge effectively generates objecthood by reconfiguring sensory mechanisms that

77 process visual inputs, thereby changing how feature space is carved up into meaningful units  
78 (Teufel & Fletcher, 2020). In other words, a given cluster of features is an object not by virtue  
79 of the features themselves but because these features are *represented as an object*. In the  
80 current study, we demonstrate that this objecthood, i.e., the fact that certain features are  
81 bound into an object representation, affects eye movements. Specifically, we show that the  
82 dynamic re-shaping of feature space by knowledge-driven perceptual organization that  
83 underlies the emergence of objecthood has a substantial influence on information sampling via  
84 eye movements in human observers.

85         The most influential early saliency models – that is, computational methods used to  
86 predict human eye-movements – largely disregarded objects, arguing that programming of  
87 eye-movements is determined by an analysis of low-level features such as luminance, colour,  
88 and orientation (Harel, Koch, & Perona, 2007; Itti & Koch, 2001). According to these early  
89 accounts, the visual system computes feature maps, which highlight areas in the image that  
90 attract fixations (Zelinsky & Bisley, 2015). Over the past 15 years, however, several studies  
91 have emphasised the importance of objects in guiding information sampling (Einhäuser, Spain,  
92 & Perona, 2008; Hayes & Henderson, 2021; Hwang et al., 2011; Nuthmann & Henderson,  
93 2010; Pajak & Nuthmann, 2013; Pilarczyk & Kuniecki, 2014; Stoll et al., 2015). For instance, in  
94 one of the early studies, Einhäuser and colleagues (2008) found that maps of object locations  
95 outperform maps derived from a low-level feature model in predicting human fixations.  
96 Moreover, human observers show a tendency to look at the centre of objects rather than their  
97 edges, contrasting with predictions from early low-level feature models (Nuthmann &  
98 Henderson, 2010; Pajak & Nuthmann, 2013; Stoll et al., 2015; Borji & Tanner, 2016; see also  
99 Vincent, Baddeley, Correani, Troscianko, & Leonards, 2009). These effects have been  
100 interpreted as demonstrations of the importance of objects in oculomotor control.

101         Other lines of evidence suggest that the fact that human observers primarily fixate at  
102 object locations can be explained by low-level mechanisms (Borji et al., 2013; Elazary & Itti,

103 2008; Kilpeläinen & Georgeson, 2018; Masciocchi et al., 2009). For instance, a recent attempt  
104 to assess the unique contribution of features vs. objects to oculomotor control suggests that  
105 object-centred effects are, at least partly, driven by low-level features that correlate with  
106 objects (Nuthmann et al., 2020). This conclusion is in line with a careful psychophysical study,  
107 suggesting that the tendency of human observers to focus on the centre of objects might be  
108 controlled by a relatively simple process that programs eye-movements towards homogeneous  
109 luminance surfaces on the basis of luminance-defined edges (Kilpeläinen & Georgeson, 2018).  
110 This result provides a potential mechanism for the finding that fixations that occur shortly after  
111 image onset tend to be located close to the stimulus centre not only for objects but also for  
112 non-objects if low-level properties are matched (Van der Linden et al., 2015). Together, these  
113 results suggest that the tendency to fixate on the centre of objects might not be related to  
114 objecthood itself but is controlled by mechanisms that respond to relatively low-level features  
115 in the input. Note, however, that the study by van der Linden and colleagues (2015) also  
116 suggests that guidance of eye-movements that are generated later after image onset might be  
117 affected by semantic aspects of object. This finding potentially indicates a time course  
118 according to which locations of early fixations are mainly determined by low-level, image-  
119 computable features while locations of later fixations might be determined by high-level object  
120 representations (see also Anderson, Ort, Kruijine, Meeter, & Donk, 2015 and Wolf & Lappe,  
121 2021).

122 Many previous studies that aim to show the contribution of objects to oculomotor  
123 control relied on a comparison of eye movements to saliency models that calculate image-  
124 computable feature maps as their null hypothesis (for example, Einhäuser et al., 2008;  
125 Pilarczyk & Kunięcki, 2014; Stoll et al., 2015). This approach has led to important insights  
126 regarding oculomotor control but is hampered by the fact that the specific methodological  
127 choices regarding the type of saliency model and object map are critical in determining the  
128 interpretation. In fact, in the previous literature, the use of different models has led to

129 categorically different conclusions, even if they have been applied to identical or very similar  
130 data sets (Borji et al., 2013; Einhauser, 2013; Einhäuser et al., 2008; Henderson et al., 2021;  
131 Henderson & Hayes, 2017; Pedziwiatr et al., 2021a, 2021b; Stoll et al., 2015). Importantly,  
132 independently of the favoured interpretation of these findings, there is a more fundamental  
133 aspect that is easily overlooked. Specifically, contrasting outputs of low-level feature models  
134 with ‘objects’, and the tendency to conceptualise these as categorically different – although  
135 possible to reconcile (Borji & Tanner, 2016; Nuthmann et al., 2020; Stoll et al., 2015) –  
136 interpretations, has concealed a fundamental similarity between these explanations. Namely,  
137 comparable to how low-level models deal with simple features, most studies implicitly treat  
138 ‘objects’ as image-computable properties. This notion is also the basis for state-of-the-art  
139 computer vision models that aim to predict human fixations (e.g., Kroner, Senden, Driessens,  
140 & Goebel, 2020; Kümmerer et al., 2017a): these models use deep convolutional neural  
141 networks trained on object recognition to extract high-level features that are directly computed  
142 from the image. In other words, the different approaches in the current eye-movement  
143 literature can be understood as lying on a continuum, with their position being defined by the  
144 type of features they emphasise. This notion is made explicit in a recent study by Schütt and  
145 colleagues (Schütt, Rothkegel, Trukenbrod, Engbert, & Wichmann, 2019): the authors explicitly  
146 conceptualised objects as high-level features that are computed in a bottom-up fashion, and  
147 contrasted their contribution to the guidance of eye-movements with the contribution of low-  
148 level features.

149 While the theoretical precision of the study by Schütt and colleagues is exceedingly  
150 helpful in clarifying the different positions, conceptualising objects as high-level features  
151 directly conflicts with current developments in object perception. Two aspects of the complex  
152 relationship between features and objects are particularly relevant: first, several recent studies  
153 demonstrate that features are not always sufficient for object representations to arise  
154 (Flounders et al., 2019; Hsieh et al., 2010; Lengyel et al., 2019, 2021; Ongchoco & Scholl,

155 2019; Teufel et al., 2018). Rather, objecthood emerges as a consequence of the interaction  
156 between current visual input and perceptual organization processes that are based on prior  
157 object-knowledge. Second, once object representations have been generated, top-down  
158 influences reconfigure the way in which even some of the earliest cortical mechanisms process  
159 low-level visual features (Christensen et al., 2015; Flounders et al., 2019; Hsieh et al., 2010;  
160 Lengyel et al., 2021, 2019; Neri, 2014, 2017; Ongchoco & Scholl, 2019; Teufel et al., 2018).  
161 For instance, psychophysical studies show that early feature-detector units are sharpened for  
162 currently relevant input based on top-down influences from object representations (Teufel et  
163 al., 2018). This reconfiguration of information processing is detectable in early retinotopic  
164 cortices (Flounders et al., 2019; Hsieh et al., 2010). Overall, these findings thus cast serious  
165 doubt on the notion that the human visual system computes image features independently of  
166 the inferred object structure of the environment (Neri, 2017).

167         This novel perspective of object perception has fundamental implications for our  
168 understanding of information sampling via eye movement. First, if objecthood emerges from the  
169 interaction between features and prior knowledge, then the question of whether objects guide  
170 eye movements cannot be answered by an approach that exclusively focuses on how image-  
171 computable feature space is carved up by the visual system, regardless of whether the  
172 considered features are low- or high-level. Second, the novel perspective of object perception  
173 means that a full understanding of the role of objects in eye-movement control has to move  
174 away from regarding feature space as static, instead taking into account the plasticity of low-  
175 level sensory processing introduced by dynamic interactions with object representations. Here  
176 we address both of these issues. We analysed gaze data from human observers viewing  
177 stimuli, which, on initial viewing, are experienced as a collection of meaningless black and white  
178 patches. After gaining relevant object knowledge, however, the observers' visual system  
179 organizes the sensory input into meaningful object representations. These stimuli allow us to  
180 test the hypothesis that eye-movements are guided by objecthood per se – i.e., the fact that



181 certain features are *represented as an object* – rather than by the high-level features associated  
182 with objects. Across three experiments (see Fig. 3 for a roadmap through them), we  
183 demonstrate that, consistent with our hypothesis, the knowledge-driven perceptual organization  
184 of identical inputs substantially re-shapes eye-movement patterns, with the selection of fixation  
185 locations being driven by a combination of image-computable features and the knowledge-  
186 dependent object representations. Moreover, these effects are already present at the first  
187 fixation. In summary, we show that a fundamental human visual behaviour – information  
188 sampling via eye movements – is guided by a dynamic interaction between image-computable  
189 features and object representations that emerge when prior object-knowledge restructures  
190 sensory input.

191

192

### Experiment 1 – Methods

193

#### Overview

194

195

196

197

198

199

200

201

202

#### Figure 1

203

*Example of a two-tone image*



204

205 *Note.* On initial viewing, this image appears as meaningless black and white patches. To be  
206 able to perceptually organize it into a meaningful percept, the reader is advised to first carefully  
207 look at the template image from which this two-tone was derived, presented on Fig. 2. An  
208 animated version of the blending between this two-tone and its template is provided in  
209 Supplementary Materials. Note that the example two-tone image is for illustration only, it was  
210 not used in the study. Image copyrights owner: author C. T.

211

212 Two-tone images provide a tool to manipulate object perception without changing the  
213 visual features of the stimulus. They are therefore ideally suited to test the hypothesis that  
214 human oculomotor control is determined by object representations that are not constituted by  
215 image-computable features but emerge via an interaction between image-computable features  
216 and prior object-knowledge. According to this idea, eye movements in response to two-tone  
217 images should be influenced by whether the observer experiences the input as an object  
218 percept. Specifically, patterns of fixations on identical two-tone images should be more similar  
219 to the ones from the corresponding template when an observer experiences the two-tone  
220 image as a meaningful object percept compared to when they experience it as meaningless

221 patches.

222 To test these predictions, we recorded eye-movements of 36 human observers who  
223 viewed two-tone images before and after being exposed to the relevant templates (Before,  
224 After, and Template conditions, respectively; see Fig. 2). In the Before condition, observers  
225 perceive two-tone images as meaningless black and white patches. In the After condition, prior  
226 object-knowledge allows them to bind patches into meaningful object percepts. Crucially, any  
227 potential differences in eye movements between the Before and the After conditions cannot be  
228 explained by image-computable features because these are identical across these conditions;  
229 the only aspect that has changed is the prior object-knowledge that observers have access to.  
230 Experiment 1 established the key effects; to exclude alternative explanations, we conducted  
231 Experiments 2 and 3 (see Fig. 3 for design details). The experiments were not preregistered.  
232 Experimental data is openly available under the following link: [link to be provided upon  
233 publication]

234

### 235 **Observers**

236 The primary units of analysis were not individual observers, but the distribution of  
237 fixations from all observers on individual images. Therefore, we selected the number of  
238 observers based on the estimation of how well our empirical fixation distributions approximate  
239 the theoretical distributions which would be obtained from the population of infinitely many  
240 observers. Previous work has shown that fixations from 18 observers provide a sufficiently  
241 good approximation for natural scenes viewed for three seconds (as in our experiment), and  
242 that further increasing the number of observers results only in marginal improvements (Judd et  
243 al., 2012). However, one of our analyses – reported in the Supplement – required splitting our  
244 sample into two groups and we therefore recruited 36 observers in total (mean age 20.06  
245 years, 7 men), ensuring sufficient amounts of data in each group after the split. All participants  
246 were Cardiff University students, had normal or corrected-to-normal vision, participated in the

247 study voluntarily, and received either money or study-credits as a reimbursement. All  
248 experiments reported in this article were approved by the Cardiff University School of  
249 Psychology Research Ethics Committee.

250

### 251 **Stimuli**

252 We used 30 pairs of images, where each pair consisted of a two-tone image and its  
253 template in greyscale. These stimuli were a subset of stimuli used in a previous study (Teufel et  
254 al., 2015), where details of template selection and two-tone image generation can be found. In  
255 brief, template images were taken from the Corel Photo library. The main objects depicted in the  
256 images were either animals (25 images), humans (three images), or animals and humans (two  
257 images). Twenty-five images depicted one main object, five images depicted two main objects.  
258 Regarding specific object-parts, seven images depicted mainly one head, two mainly two heads,  
259 18 depicted a head with a full body, and three images depicted two full bodies with heads. Two-  
260 tones were generated by smoothing and binarising template images. A good two-tone image  
261 should be perceived as a collection of meaningless patches prior to seeing its template but  
262 observers should be able to easily bind the stimulus into a coherent percept of an object after  
263 they see the template. Extensive tests on naïve observers were conducted to select both the  
264 template images, and the parameters of smoothing and binarisation that guarantee that the  
265 created two-tones have these desired properties. Note that two-tone images are different from  
266 Mooney stimuli (Mooney 1957). By contrast to two-tone images, Mooney stimuli can be, and are  
267 designed to be, recognized spontaneously (without need for prior knowledge).

268

### 269 **Experimental setup**

270 The experiment was conducted in a dark testing room. Participants sat 56 cm from the  
271 monitor, with their head supported by a chin and forehead rest. Their eye-movements were  
272 recorded using an EyeLink 1000+ eye-tracker (with 500 Hz sampling rate) placed on a tower

273 mount. The experiment was controlled by in-house developed code written in Matlab R2016b  
274 (Mathworks, Natick, MA) and using the Psychophysics Toolbox Version 3 (Brainard, 1997;  
275 Kleiner et al., 2007). Images were presented centrally on the screen, against a mid-grey  
276 background. Images measured 21.9 degrees of visual angle (788 pixels) horizontally and 14.6  
277 degrees (526 pixels) vertically.

278

## 279 **Procedure**

280         The experiment consisted of ten blocks; a single block is schematically illustrated in  
281 Fig. 2. Before the start of the procedure, a 13-point eye-tracker calibration and validation was  
282 conducted. Each block started with the Before condition, in which three two-tones were  
283 presented in a sequence, each for 3 seconds. Observers were instructed to carefully look at  
284 these images, but they were not specifically told to search for objects. Two-tone images were  
285 preceded by a centrally-located fixation-dot displayed for 1 second. They were followed by a  
286 visual analogue scale, which observers adjusted by pressing 'z' and 'm' buttons on a keyboard  
287 to indicate how meaningful they experienced the two-tone image to be. The instruction given to  
288 the observers prior to the experiment was also displayed above the scale, saying: 'Please  
289 indicate how clearly the scene or object in the image appeared to be.' The scale was  
290 continuous, with the following labels placed at five linearly spaced points above the scale:  
291 'Very unclear', 'Unclear', 'Neither clear nor unclear', 'Clear', 'Very clear'. Meaningfulness ratings  
292 were used as a manipulation check. After each rating, a blank screen was displayed for 500  
293 ms. The Before condition was followed by the Template condition, in which template images  
294 were displayed while eye-movements were recorded – again, each for 3 seconds, preceded by  
295 a fixation dot. After the Template condition, we ensured that observers had enough object-  
296 knowledge to bind two-tone images into meaningful object percepts by presenting six cycles of  
297 dynamic blending between two-tones and their templates (Blending Phase). Each cycle began  
298 with the presentation of a template image for two seconds. This was then linearly blended into

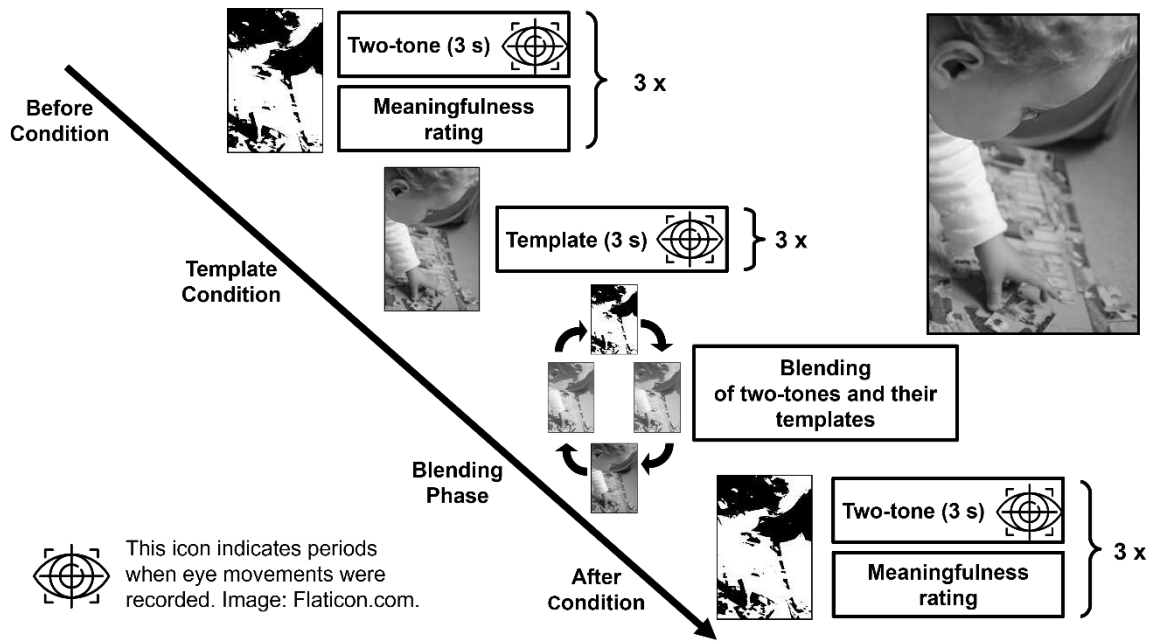
299 the corresponding two-tone image, with the full transition from template to two-tone taking 4  
300 seconds. The two-tone image remained on the screen for 2 seconds and then was blended  
301 back into the template, remaining on the screen for another 2 seconds. Each of the three  
302 image-pairs used in a block was presented in a full blending procedure twice with the order  
303 pseudorandomised such that the same pair was never used twice in a row. The subsequent  
304 cycles of blending were separated with a blank screen presented for 500 ms. After the  
305 Blending Phase, the After condition was presented, which was identical to the Before condition  
306 except that images were presented in a newly randomized order. There was a break every two  
307 blocks, and the eye-tracker was re-calibrated. For each observer, images were assigned to  
308 blocks randomly and were presented in a pseudo-random order within each block. The  
309 pseudo-randomization ensured that the image shown last in the Blending Phase was never  
310 presented at the beginning of the After condition. Total experiment time was ~50 minutes.

311 Instructions were delivered verbally and on-screen. Key elements of the procedure were  
312 illustrated visually: observers were shown a single two-tone image (which was not used in the  
313 actual experiment), rated its meaningfulness, viewed the blending procedure with the template  
314 and, finally, viewed the same two-tone again and were asked to provide a meaningfulness  
315 rating.

316

## 317 **Figure 2**

318 *Experiment 1 – Outline of a single experimental block*



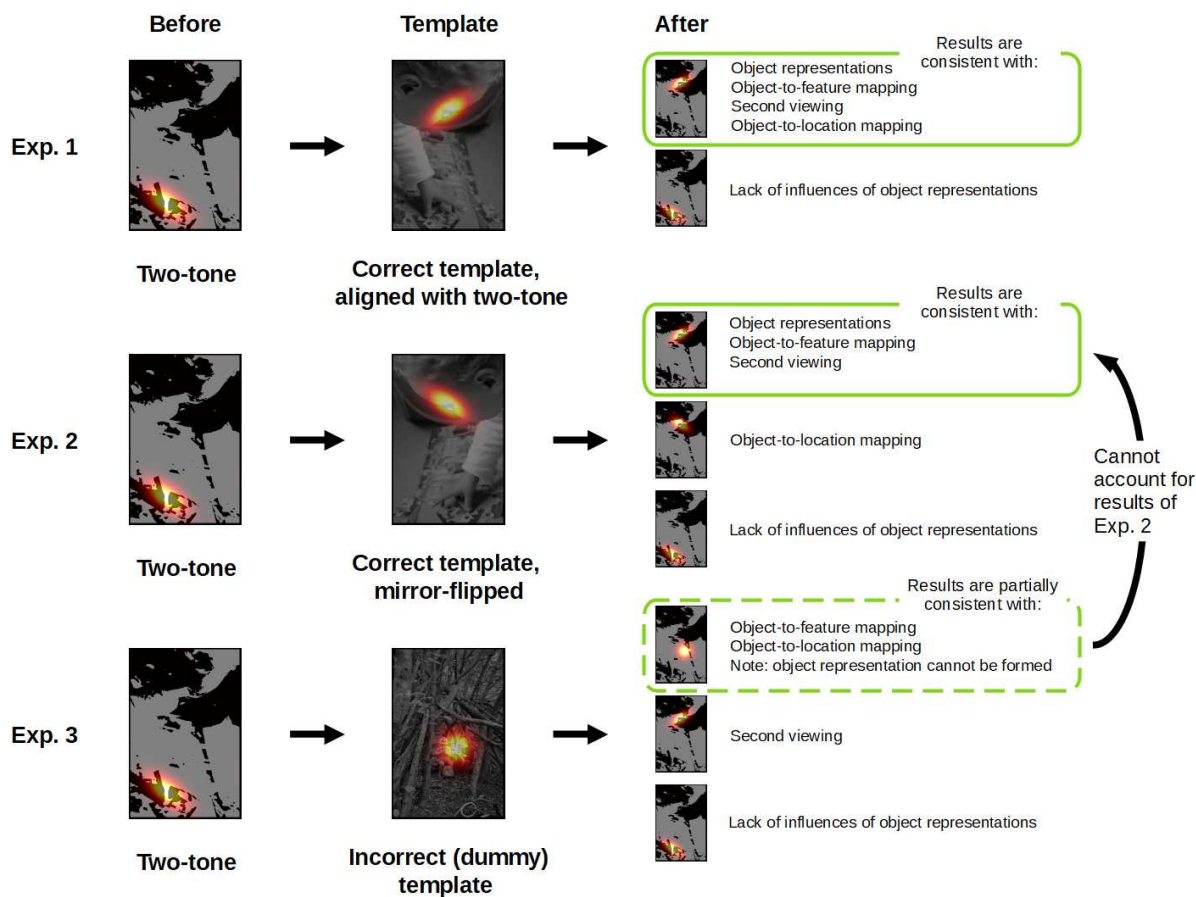
319

320 *Note.* In each block, observers first free-viewed three two-tone images (Before condition). After  
 321 presentation of each image, they were asked to rate its perceived meaningfulness. Then, the  
 322 grayscale templates of these three two-tones were presented (Template condition). In the next  
 323 part of the block, observers viewed the two-tones gradually blended with their templates six  
 324 times (Blending Phase). The After condition was identical to the Before condition in all aspects  
 325 except for the order of presentation of the two-tone images. In the upper right corner, the  
 326 template of the two-tone image from Fig. 1 is presented (copyrights owner: author C. T.).

327

328 **Figure 3**

329 *Summary of key experimental manipulations, predictions, and findings of Experiment 1, 2, and 3*



330

331 **Note.** The heatmaps superimposed over the example stimuli illustrate hypotheses we test in

332 our experiments. The After column illustrates potential experimental outcomes, with green

333 rectangles indicating interpretations consistent with the results for each experiment. All three

334 experiments had identical designs except for the type of image shown in the Template

335 condition and in the Blending Phase. In Experiment 1, the original grayscale photograph used

336 to generate the two-tone image provided observers with the prior object-knowledge required to

337 organise the two-tone image into a coherent object percept in the After condition. We found

338 that gaze guidance in the After condition was similar to that in the Template condition (first row,

339 right top panel), suggesting that knowledge-driven perceptual organization is an important

340 driver of oculomotor control. In Experiments 2 and 3, we excluded potential alternative

341 explanations. In Experiment 2, we presented mirror-flipped template images. This manipulation

342 allowed us to exclude the possibility that when viewing the templates, observers learned the



343 position of objects in the images, and re-visited these locations in the After condition. In  
344 Experiment 3, ‘dummy templates’ unrelated to the two-tone images were presented, which  
345 allowed us to exclude the possibility that second-viewing of the two-tone images could explain  
346 the results. Moreover, this design allowed us to test whether observers had learned to map the  
347 features of a two-tone image to locations of objects in the template images. We found a small  
348 effect consistent with this idea, but it was too small to fully account for the main findings.

349

### 350 **Data pre-processing and analysis methods**

351 The default EyeLink algorithm was used to extract fixation-locations from the eye-  
352 movement recordings. Further data pre-processing was done in Matlab. For each image, we  
353 discarded the initial fixation that was directed at the fixation-dot presented before image onset.  
354 We also discarded fixations not landing within the image-boundaries. Further details regarding  
355 data exclusions can be found in the Data exclusion section of the Supplement. For each image  
356 in each condition, we generated heatmaps (see examples on Fig. 5E) by smoothing the  
357 discrete distribution of fixations with a Gaussian filter, cutoff frequency of  $-6\text{dB}$   
358 (implementation provided by Bylinskii and colleagues; Kümmerer et al., 2020), and then  
359 normalizing the smoothed distribution to the zero-one range.

360 The majority of our analyses focused on the similarity between two heatmaps. As a  
361 similarity index, we calculated Pearson’s linear correlation coefficient using Matlab  
362 implementation (Kümmerer et al., 2020). This measure is intuitive, commonly used in the  
363 literature (Wilming, Betz, Kietzmann, & König, 2011), and its values have a straightforward  
364 interpretation. In the current study, values ranged between zero and one, with one indicating  
365 that two heatmaps are identical and zero indicating a maximal dissimilarity. In the Supplement,  
366 we provide the results of key analyses using a different metric to quantify the similarity  
367 between two heatmaps, the histogram intersection (SIM; Bylinskii, Judd, Oliva, Torralba, &  
368 Durand, 2016), showing a similar pattern of results. For statistical comparisons, we primarily

369 relied on standard null-hypothesis-significance-testing techniques implemented in R (R Core  
370 Team, 2020) and Matlab. Unless otherwise stated, the t-tests reported throughout the text are  
371 paired-sample t-tests. In order to assess the amount of evidence for a lack of a difference  
372 between groups of measurements, we used Bayes factors (BFs) calculated using bayesFactor  
373 R package (Morey & Rouder, 2018).

374

## 375 **Experiment 1 – Results**

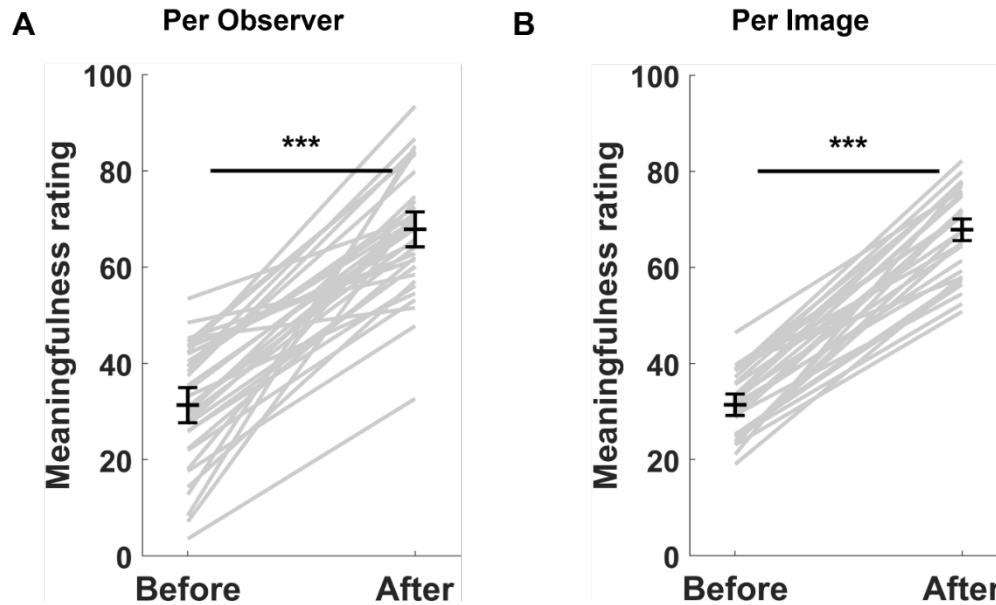
### 376 **Manipulation check: analysis of meaningfulness ratings**

377 In the Before and After conditions, observers rated the perceived meaningfulness of two-  
378 tone images. Averaging these ratings per image showed that the two-tones were perceived as  
379 more meaningful in the After compared to the Before condition (Fig. 4A and B;  $t(29) = 23.84$ ,  $p <$   
380  $0.001$ ; mean difference  $M_{diff} = 0.36$ , 95% confidence interval  $CI = [0.33, 0.4]$ ). The same pattern  
381 of results held when the ratings were averaged per observer ( $t(35) = 14.42$ ,  $p < 0.001$ ;  $M_{diff} =$   
382  $0.37$ , 95%  $CI = [0.31, 0.42]$ ). These results provide a manipulation check, suggesting that  
383 observers are able to organize two-tone images into meaningful object representations after but  
384 not before acquiring relevant prior object-knowledge.

385

### 386 **Figure 4**

387 *Meaningfulness ratings for two-tone images in the Before and After conditions averaged per*  
388 *observer (A) and per image (B)*



389

390 *Note.* The following conventions are used in this and all remaining figures: asterisks on plots391 indicate p-values: \*\*\* indicates  $p \leq 0.001$ , \*\* indicates  $p \leq 0.01$ , \* indicates  $p \leq 0.05$ , and 'n.s.'

392 indicates the lack of statistical significance. Grey lines indicate values for individual observers

393 (panel A) and images (panel B). Black horizontal bars indicate means. They are surrounded

394 with 95% confidence intervals for within-subjects designs, calculated using Cousineau-Morey

395 method (Cousineau, 2005; Morey, 2008).

396

397 **Analysis of similarity between heatmaps**

398 If knowledge-dependent object representations drive eye movements, the spatial

399 distribution of fixations recorded in response to two-tone and template images should be more

400 similar when two-tone images elicit object representations (After condition) compared to when

401 they do not (Before condition). To test this hypothesis, we compared the similarities of

402 heatmaps across pairs of conditions (Fig. 5A). As predicted, we found a higher similarity

403 between the Template-After pair ( $M = 0.90$ ,  $SD = 0.07$ ) compared to the Template-Before pair404 ( $M = 0.72$ ,  $SD = 0.13$ ;  $t(29) = 8.39$ ,  $p < 0.001$ ; mean difference  $M_{diff} = 0.18$ , 95% CI = [0.14,

405 0.22]). This result suggests that gaze patterns in response to two-tone images more closely

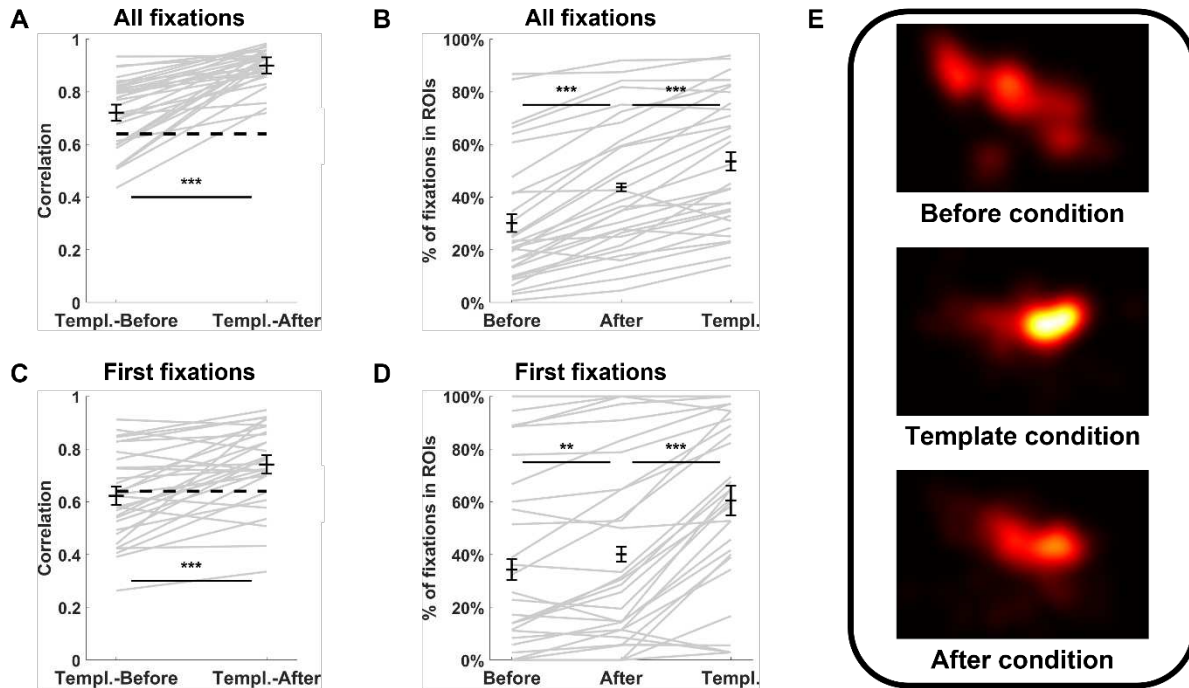
406 resemble eye movements from the templates when the two-tones were perceived as  
407 containing meaningful objects, as compared to when they were perceived as meaningless  
408 patches.

409         While there was a clear difference in similarity between the two pairs, at first glance the  
410 Template-Before similarity might seem unexpectedly high. Importantly, however, the distribution  
411 of fixations on images is not only determined by the characteristics of the visual input but also  
412 by general factors that are independent of the image (Tatler & Vincent, 2009). One key general  
413 factor is the centre bias, a tendency of humans to look at the centre of an image rather than  
414 regions closer to the edges (Tatler, 2007). A meaningful evaluation of the difference in  
415 similarities between Template-Before and Template-After pairs therefore requires a baseline  
416 that accounts for this bias. Given that there is no consensus on exactly how to model centre  
417 bias (Hayes & Henderson, 2020), and that systematic studies of centre bias only exist for a  
418 limited number of combinations of image sizes and aspect ratios (Clarke & Tatler, 2014), we  
419 adopted a data-driven approach to derive a centre bias. Specifically, we modelled a centre bias  
420 for our data by creating a single heatmap (labelled 'Centre') from all fixations registered  
421 throughout the experiment. The rationale for this approach is that by averaging across all  
422 images and all observers, the remaining heatmap should include only those factors that are  
423 general to all images and observers (i.e., centre bias) in our dataset. We found a statistically  
424 robust difference in similarity scores between the Template-Centre and Template-Before pairs  
425 (Template-Centre:  $M = 0.64$ ,  $SD = 0.16$ ; Template-Before:  $M = 0.72$ ,  $SD = 0.13$ ;  $t(29) = 2.40$ ,  $p =$   
426  $0.023$ ;  $M_{diff} = 0.08$ ,  $95\% CI = [0.01, 0.14]$ ). Importantly, however, this difference was small,  
427 suggesting that a centre bias explained most, but not all, of the Template-Before similarity.

428

## 429 **Figure 5**

430 *Results of Experiment 1*



431

432 *Note.* A) Similarities between heatmaps from template and two-tone images, where the two-

433 tones were viewed either in the Before or in the After condition. The dashed horizontal line

434 illustrates the baseline, i.e., the expected similarity with the Template condition based purely on

435 centre bias. B) Proportion of fixations landing within the regions-of-interest (ROIs) in each

436 condition. ROIs included important object parts (e.g., the heads of depicted animals). C, D) The

437 same analyses as on panels A and B but conducted including only first fixations from the Before

438 and After conditions. E) Sample heatmaps illustrating the distributions of fixations in all three

439 conditions of Experiment 1 for one two-tone/template pair. These maps were created from all

440 fixations registered on the images. Pixel values of all three maps were jointly normalised to

441 zero-one range, so colour values (indicating fixation densities) are comparable across panels.

442

443 We ran a further analysis (full details in Supplement) to address the influence of knowledge-

444 dependent object representations by comparing heatmaps from identical visual inputs only. In

445 other words, instead of analysing the similarity between heatmaps from a two-tone image and

446 its template image (different visual inputs), we evaluated the similarities in heatmaps when the

447 same two-tone image was viewed in the Before and the After conditions (identical visual inputs).  
448 The findings provide further support for the influence of object-knowledge on gaze guidance  
449 (see supplement for details).

450

### 451 **Regions-of-interest (ROI) analysis**

452         The analyses of heatmap similarities suggests that prior object-knowledge contributes to  
453 eye-movement control. We used a region-of-interest (ROI) analysis to assess in a more fine-  
454 grained manner the extent to which changes in fixation patterns related directly to object  
455 representations. We exploited the fact that animal and human heads are known to attract  
456 fixations in natural scenes (Cerf, Paxon Frady, & Koch, 2009; Drewes, Trommershäuser, &  
457 Gegenfurtner, 2011). On each template, we manually labelled each pixel associated with a head  
458 (recall that all templates depicted animals and/or humans). The resulting masks, which covered  
459 9% of the image area on average (SD = 12%, median = 3%), served as the ROIs for the  
460 template and its associated two-tone image. The average distance of the centre of gravity of  
461 each masks (as determined by Matlab function *regionprops*) to the image centre was 3.83  
462 degrees of visual angle (SD = 2.91) and the distance to the central vertical image axis was 2.19  
463 degrees (SD = 2.23). For each image and condition, we calculated the proportion of fixations  
464 landing within the ROIs (Fig. 5B). This metric increased in the After compared to the Before  
465 condition, indicating that changes in fixations were object-specific (Before: M = 30%, SD = 24;  
466 After: M = 44%, SD = 25;  $t(29) = 8.64$ ,  $p < 0.001$ ;  $M_{diff} = 0.14$ , 95% CI = [0.1, 0.17]).  
467 Furthermore, there were more fixations within the ROIs in the Template compared to the After  
468 condition (Template: M = 54%, SD = 25;  $t(29) = 6.02$ ,  $p < 0.001$ ;  $M_{diff} = 0.1$ , 95% CI = [0.06,  
469 0.13]). Overall, the ROI analysis provides clear evidence to suggest that the influence of  
470 knowledge-dependent object representations on fixation patterns is object-specific.

471

**472 Analysis of first fixations**

473 In order to assess the time-course of the influence of knowledge-dependent object  
474 representations on oculomotor control, we repeated our previous analyses exclusively for first  
475 fixations. This restriction did not change the overall pattern of the results (see Fig. 5C and D),  
476 suggesting that even first fixations were influenced by object representations that emerged as a  
477 consequence of the observer's prior knowledge. Specifically, the statistical analysis showed that  
478 for first fixations, the similarity between Template and After was higher than for Template and  
479 Before (Template-After:  $M = 0.74$ ,  $SD = 0.15$ ; Template-Before:  $M = 0.62$ ,  $SD = 0.17$ ;  $t(29) =$   
480  $4.91$ ,  $p < .001$ ;  $M_{diff} = 0.12$ ,  $95\% CI = [0.07, 0.17]$ ). This finding was corroborated by an ROI  
481 analysis of first fixations: the proportion of first fixations landing on ROIs was higher in the After  
482 than in the Before condition, and also higher in Template than in After (Before:  $M = 34\%$ ,  $SD =$   
483  $34$ ; After:  $M = 40\%$ ,  $SD = 35$ ; Template:  $M = 60\%$ ,  $SD = 32$ ; Before-After:  $t(29) = 3.61$ ,  $p =$   
484  $0.001$ ;  $M_{diff} = 0.06$ ,  $95\% CI = [0.03, 0.09]$ ; Template-After:  $t(29) = 6.41$ ,  $p < 0.001$ ;  $M_{diff} = 0.2$ ,  
485  $95\% CI = [0.14, 0.27]$ ). Taken together, these results suggest that knowledge-dependent object  
486 representations emerge fast enough to influence even the first eye-movements after stimulus  
487 onset.

488

**489 Analysis of combined effects of image-computable features and prior knowledge**

490 Our analyses so far indicate that knowledge-dependent object representations play a  
491 role in gaze guidance, beginning with the first fixation after image onset. However, these  
492 analyses do not assess the role of the interaction between image-computable features and  
493 object representations. In order to address this point, we capitalized on common and distinct  
494 characteristics shared between the After condition and each of the remaining conditions (Before  
495 and Template). In particular, image-computable features of Before and After conditions are  
496 identical, but they differ in the extent to which observers experienced object representations.  
497 Specific similarities in fixation patterns between Before and the After conditions, which go

498 beyond general factors such as centre bias, can therefore be attributed to the image-  
499 computable features of two-tone images. Conversely, the After and the Template conditions  
500 have the reverse relationship: they lead to similar object representations but differ in image-  
501 computable features. We exploited this situation to characterize the contribution of these gaze-  
502 guidance factors in the After condition.

503 For this purpose, we created linear combinations of heatmaps from the Before and  
504 Template conditions to compare with the heatmaps of the After condition (Fig. 6). Each new  
505 linear-combination heatmap was calculated from the Before and the Template conditions'  
506 heatmaps, using the formula:

507

$$508 \quad w_{Template} * heatmap_{Template} + w_{Before} * heatmap_{Before} \quad (1)$$

509

510 where  $w$  is a weight for the heatmap indicated by the subscript. Incorporating the normalization  
511 assumption ( $w_{Template} + w_{Before} = 1$ ), we created a continuum of heatmaps spanning the range  
512 between being fully determined by the Template heatmap to being fully determined by the  
513 Before heatmap. This continuum was uniformly sampled with a step-size of 0.05. This  
514 procedure led to a set of heatmaps, which capture factors driving eye movements in the Before  
515 and the Template conditions to varying degrees. Evaluating the similarity of these new  
516 heatmaps with those from the After condition allowed us to determine the relative contribution of  
517 image-computable features and object representations to gaze guidance in the After condition.  
518 To focus on the time course, we conducted this analysis separately for first fixations and for all  
519 the remaining fixations.

520 The results of this similarity analysis suggest that both first and all remaining fixations in  
521 the After condition were guided synergistically by image-computable features and object  
522 representations (Fig. 7). The linear-combination heatmaps that had the highest similarity with  
523 the first fixations in the After condition showed an influence from the Template heatmap but also



524 had a substantial contribution from the Before heatmap ( $w_{\text{Template}} = 0.4$ ,  $w_{\text{Before}} = 0.6$ ; mean  
525 correlation  $M = 0.85$ ,  $SD = 0.06$ ; see Fig. 7A). Statistical analyses indicated that the heatmaps  
526 from the After condition were more similar to this optimal linear-combination heatmap than to  
527 either the Before or the Template conditions alone (Optimal-After vs. Before-After:  $t(29) = -2.67$ ,  
528  $p = 0.012$ ;  $M_{\text{diff}} = 0.03$ , 95% CI = [0.01, 0.04]; Optimal-After vs. Template-After:  $t(29) = 5.70$ ,  $p$   
529  $< 0.001$ ;  $M_{\text{diff}} = 0.11$ , 95% CI = [0.07, 0.15]).

530         The findings for all remaining fixations from the After condition were similar (Fig. 7B).  
531 However, the linear combinations that were optimal for these fixations were more strongly  
532 influenced by the Template heatmap ( $w_{\text{Template}} = 0.65$ ,  $w_{\text{Before}} = 0.35$ ; mean correlation  $M =$   
533  $0.95$ ,  $SD = 0.03$ ). Yet, even for these later fixations, there was a substantial influence of image-  
534 computable factors as captured by the Before heatmaps. This idea is supported by the  
535 statistical analysis, which indicates that the heatmaps from the After condition were more similar  
536 to the optimally combined heatmaps compared to both the Before and the Template condition  
537 alone (Optimal-After vs. Before-After:  $t(29) = 6.49$ ,  $p < 0.001$ ;  $M_{\text{diff}} = 0.09$ , 95% CI = [0.06,  
538 0.12]; Optimal-After vs. Template-After:  $t(29) = 5.48$ ,  $p < 0.001$ ;  $M_{\text{diff}} = 0.05$ , 95% CI = [0.03,  
539 0.06]).

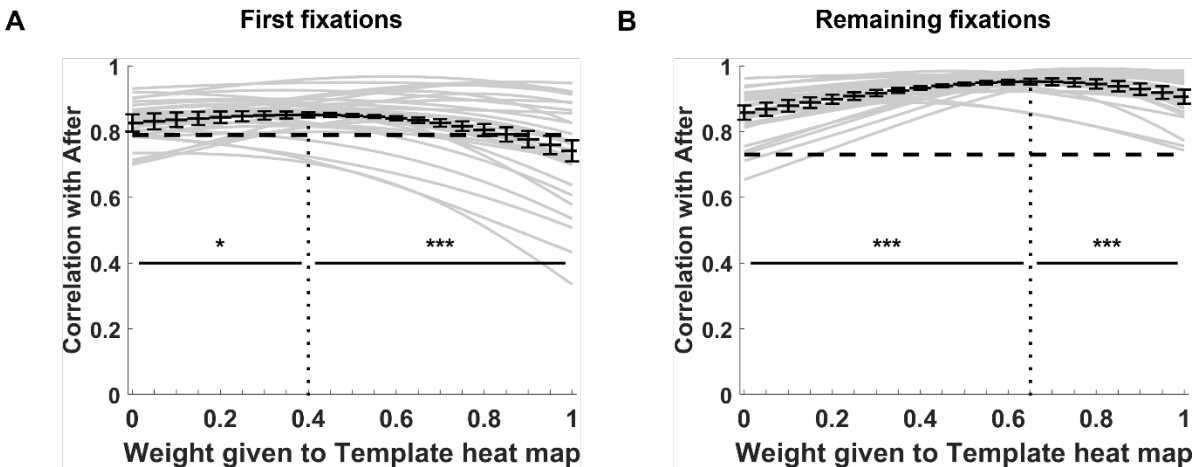
540         Overall, the analysis suggests that image-computable features and object  
541 representations guide eye movements in a synergistic manner (see also Borji & Tanner, 2016).  
542 The contribution of these two factors vary over time, with object representations playing a less  
543 important role in first fixations than in later fixations. Yet, both factors already influence first  
544 fixations.

545

## 546 **Figure 6**

547 *Linear combination analysis – illustration for a single two-tone image*





561

562 *Note.* A) Similarities obtained when only first fixations from the After condition are considered.

563 B) The same analysis but for all the remaining fixations (i.e., without first) from the After

564 condition. The weights of the linear combinations for which the similarity is maximal are

565 indicated by the dotted vertical lines. Dashed vertical lines on both panels indicate the baseline,

566 i.e., the average similarities of the respective After heatmaps to centre bias model ( $M = 0.79$ ,567  $SD = 0.09$  for first fixations;  $M = 0.73$ ,  $SD = 0.14$  for the remaining ones).

568

569 **Analysis of other characteristics of oculomotor behaviour**

570 In our final analyses of Experiment 1, we assessed the extent to which knowledge-

571 dependent object representations affect characteristics of eye movements that might be

572 indicative of a more fundamental change in the observers' information-sampling strategy. First,

573 we calculated the mean number of fixations, average fixation duration (in seconds), and

574 average Euclidean distance between consecutive fixations (interfixation distance, in degrees of

575 visual angle) per image, and compared them across conditions (Fig. 8). Compared to the Before

576 condition, the After condition showed a decrease in the number of fixations (values summed

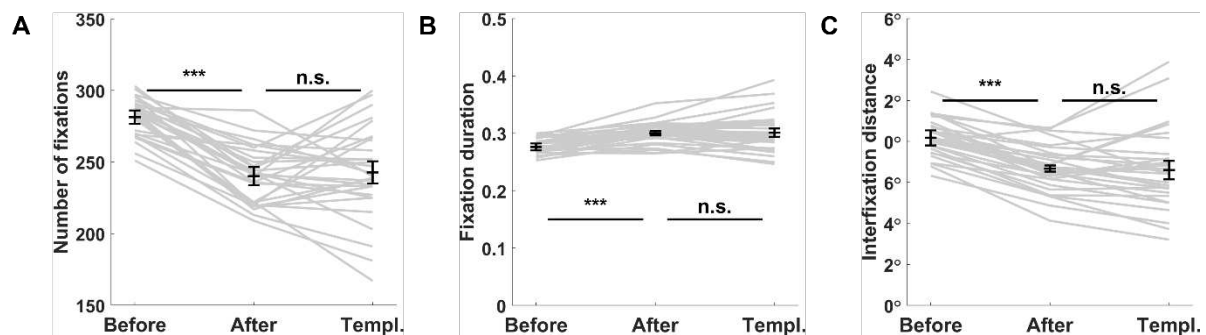
577 across observers separately for each image; Before:  $M = 281.37$ ,  $SD = 13.22$ ; After:  $M = 240.10$ ,578  $SD = 19.32$ ;  $t(29) = 12.76$ ,  $p < 0.001$ ;  $M_{diff} = 41.27$ ,  $95\% CI = [34.65, 47.88]$ ), an increase in the579 fixation duration (Before:  $M = 0.28$ ,  $SD = 0.01$ ; After:  $M = 0.30$ ,  $SD = 0.02$ ;  $t(29) = -8.22$ ,  $p <$

580 0.001;  $M_{diff} = -0.02$ , 95% CI = [0.02, 0.03]), and a decrease in interfixation distance (Before:  $M$   
 581 = 4.09,  $SD = 0.45$ ; After:  $M = 3.34$ ,  $SD = 0.55$ ;  $t(29) = 11.24$ ,  $p < 0.001$ ;  $M_{diff} = 0.75$ , 95% CI =  
 582 [0.61, 0.89]). We did not find statistically significant differences between the Template and the  
 583 After conditions for any of these metrics (number of fixations:  $t(29) = -0.50$ ,  $p = 0.621$ ;  $M_{diff} = -$   
 584 2.67, 95% CI = [-13.58, 8.25]; fixation duration:  $t(29) = -0.24$ ,  $p = 0.816$ ;  $M_{diff} = 0$ , 95% CI = [-  
 585 0.01, 0.01]; interfixation distance:  $t(29) = 0.32$ ,  $p = 0.755$ ;  $M_{diff} = 0.04$ , 95% CI = [-0.19, 0.27];  
 586 descriptive statistics for these three respective characteristics for Template condition:  $M =$   
 587 242.77,  $SD = 31.76$ ;  $M = 0.30$ ,  $SD = 0.03$ ;  $M = 3.3$ ,  $SD = 0.96$ ).

588

589 **Figure 8**

590 *Number of fixations (A), fixation duration (B), and interfixation distance measured in degrees of*  
 591 *a visual angle (C)*



592

593 *Note.* All three were calculated per image and compared between conditions.

594

595 These findings are consistent with the idea that observers shift from exploring the whole  
 596 stimulus in the Before condition towards extracting information only from selected parts in the  
 597 After and Template conditions. To further substantiate this interpretation, we calculated the  
 598 normalized entropy for the heatmaps in the different conditions (Fig. 9A). This measure is  
 599 thought to index the extent to which an observer's behaviour is exploratory (Gameiro, Kaspar,  
 600 König, Nordholt, & König, 2017; Kaspar et al., 2013). Normalized entropy was lowest in the

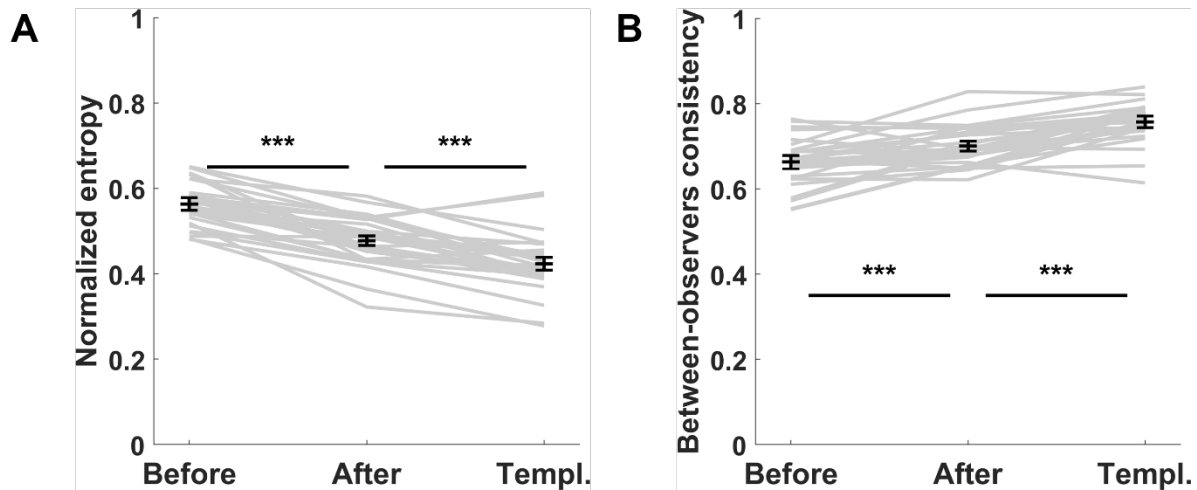
601 Template condition, increased in the After condition, and was highest in the Before condition  
602 (Before:  $M = 0.56$ ,  $SD = 0.05$ ; After:  $M = 0.48$ ,  $SD = 0.06$ ; Template:  $M = 0.42$ ,  $SD = 0.07$ ;  
603 Before-After:  $t(29) = 9.92$ ,  $p < 0.001$ ;  $M_{diff} = 0.09$ , 95% CI = [0.07, 0.10]; After-Template:  $t(29) =$   
604 6.28,  $p < 0.001$ ;  $M_{diff} = 0.05$ , 95% CI = [0.04, 0.07]). In other words, observers showed the  
605 highest exploratory behaviour in the Before condition, followed by the After and the Template  
606 condition.

607 In our final analysis, we wanted to know if object representations would result in more  
608 homogenous gaze behaviour across observers (Fig. 9B). We quantified between-observers  
609 consistency by averaging the similarity between each observer's individual heatmap to the  
610 heatmaps of all remaining observers (Lyu et al., 2020). This metric increased both between the  
611 Before and After conditions and between the After and Template conditions (Before:  $M = 0.66$ ,  
612  $SD = 0.05$ ; After:  $M = 0.7$ ,  $SD = 0.05$ ; Template:  $M = 0.76$ ,  $SD = 0.05$ ; Before-After:  $t(29) = 3.96$ ,  
613  $p < 0.001$ ;  $M_{diff} = 0.04$ , 95% CI = [0.02, 0.06]; After-Template  $t(29) = 6.96$ ,  $p < 0.001$ ;  $M_{diff} =$   
614 0.06, 95% CI = [0.04, 0.07]), suggesting that object representations increase consistency in  
615 information-sampling behaviour across observers.

616

## 617 **Figure 9**

618 *Normalized entropy and between-observers consistency*



619

620 *Note.* A) Normalized entropy of fixation distributions (in arbitrary units) as a measure of their  
 621 spread. Higher values indicate more exploratory behaviour of observers. B) Between-observers  
 622 consistency in selecting fixation targets measured by how similar (on average) fixations of a  
 623 single observer were to fixations of all the remaining observers pooled together.

624

625

### Experiment 1 – Discussion

626 In Experiment 1, we measured eye-movements in response to grayscale images of scenes  
 627 containing objects and two-tone images derived from these templates. On initial viewing, two-  
 628 tone images are experienced as meaningless black-and-white patches. Once an observer has  
 629 acquired relevant prior object-knowledge, however, the visual system organizes the patches into  
 630 a coherent percept of an object. We demonstrate that, when a two-tone image is perceived as  
 631 showing a coherent object rather than meaningless patches, gaze guidance changes in several  
 632 ways. First, and most importantly, fixation patterns on two-tone images become more similar to  
 633 those measured in response to the template when two-tones lead to object representations vs.  
 634 when they are experienced as meaningless patches. Moreover, fixation locations become more  
 635 object-specific. Importantly, however, we also demonstrate that object representations do not  
 636 fully dominate gaze guidance, but that image-computable feature space and object  
 637 representations interact in determining where people look. While the data suggest a specific

638 temporal development of this interaction, we also observe that the influence of knowledge-  
639 dependent object representations is already present in the first eye-movement after image  
640 onset, suggesting that the emergence of knowledge-driven object representations precedes the  
641 first eye-movement. Object representations also lead to fewer fixations, longer fixation  
642 durations, shorter interfixation distances as well as a less exploratory pattern of eye movements  
643 and more consistency across observers. Overall, these results suggest that object  
644 representations, which are not fully determined by image-computable features but depend on an  
645 observer's prior object-knowledge have a substantial influence on eye movements. Note that  
646 the images were presented in batches of three (see the *Procedure* section), ensuring that they  
647 were not fully predictable. These results are therefore unlikely to be explained by planning of  
648 eye movements done before the onset of the image in the After condition.

649         It is, however, possible that the change in fixation patterns observed in Experiment 1  
650 were caused by a memory process unrelated to knowledge-driven perceptual organization.  
651 Specifically, it has been suggested that eye movements performed during memory retrieval of  
652 an image resemble the eye movements performed when seeing this stimulus for the first time  
653 (Noton & Stark, 1971; see Wynn, Shen, & Ryan, 2019 for a recent review and Foulsham &  
654 Kingstone, 2013 for criticism). According to this alternative explanation, two-tone images in the  
655 After condition might have acted as cues that triggered the retrieval of the corresponding  
656 template, and this retrieval might have been accompanied by the re-enactment of gaze  
657 behaviour from the Template condition. A simpler but overall similar alternative explanation of  
658 the results from Experiment 1 might suggest that memory-retrieval of template images resulted  
659 in the observers voluntarily directing their gaze towards locations in the two-tone images, which  
660 they remembered to be occupied by objects. According to both explanations, the factor driving  
661 changes in eye movements in the After condition is the mapping of objects to locations that the  
662 observers remember from the Template condition, rather than perceptual organization induced  
663 by prior object-knowledge. To exclude these alternative explanations, which we label the

664 'object-to-location mapping' interpretation, we conducted Experiment 2.

665

666

## Experiment 2

### 667 Overview

668 Experiment 2 was identical to Experiment 1 in all aspect except that the template images  
669 were flipped along the vertical axis ('mirror-flipped') from left to right. Consequently, the screen  
670 locations occupied by objects differed between the Template condition and the remaining  
671 conditions. This simple manipulation allowed us to adjudicate between the different alternative  
672 interpretations mentioned in the previous section: according to the object-to-location mapping  
673 hypothesis, which suggests that observers merely revisited the parts of the display, which  
674 contained objects during the presentation of template images, we would expect a high  
675 similarity between heatmaps from the After and Template conditions, despite the lack of  
676 overlap in spatial location of objects in these two conditions. If, however, the effects observed  
677 in Experiment 1 were attributable to knowledge-dependent object representations, we would  
678 expect the similarity between the After and Template conditions to be low (see Fig. 3 for  
679 illustration). Moreover, by mirror-flipping the heatmaps obtained from the mirror-flipped  
680 templates, we would expect an increase in similarity to levels seen in Experiment 1 (because  
681 this leads to a re-alignment of heatmaps from templates and two-tones).

682

### 683 Experiment 2 – Method

684 A separate set of 18 Cardiff University students (mean age 19.5 years, 5 men), who did  
685 not participate in Experiment 1, served as observers. The design of Experiment 2 was identical  
686 to that of Experiment 1 except that the template images were flipped along the vertical axis from  
687 left to right for all parts of the experiment. Additionally, during the Blending Phase, the two-tones  
688 were flipped such that two-tones and templates were aligned. This condition is labelled  
689 FlippedTemplate. Observers were not explicitly informed about the flipping; the instructions



690 were identical to those in Experiment 1.

691

692

## Experiment 2 – Results

693

### Controlling for the effects of object-to-location mapping

694

695

696

697

698

699

Similar to Experiment 1, the meaningfulness ratings provided by the observers after viewing each two-tone were higher in the After condition than the Before condition both when we averaged them per observer ( $t(17) = 6.62, p < 0.001; M_{diff} = 0.24, 95\% CI = [0.16, 0.31]$ ) and per image ( $t(29) = 16.74, p < 0.001; M_{diff} = 0.24, 95\% CI = [0.21, 0.27]$ ). This result indicates that observers were able to bind the two-tone images into meaningful percepts despite viewing templates, which were presented in a mirror-flipped manner.

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

The results of the eye-movements data analysis were inconsistent with the object-to-location hypothesis but provided support for the idea that knowledge-dependent object representations influence eye movements (see Fig. 10). In particular, by contrast to the analogous analysis in Experiment 1, heatmap similarities did not differ when comparing the FlippedTemplate-Before pair vs. the FlippedTemplate-After pair (FlippedTemplate-Before:  $M = 0.46, SD = 0.22$ ; FlippedTemplate-After:  $M = 0.48, SD = 0.22$ ;  $t(29) = 1.45, p = 0.158$ ;  $M_{diff} = 0.03, 95\% CI = [-0.01, 0.06]$ ). A BF of 0.50 suggested that the data provided evidence in favour of there being no difference between conditions, but that this evidence was weak. Importantly, once the heatmaps from the template and two-tone images were re-aligned, by flipping the heatmaps of the FlippedTemplate condition, the similarity between the RealignedTemplate and the After condition was higher than the similarity between RealignedTemplate and Before (RealignedTemplate-Before:  $M = 0.68, SD = 0.15$ ; RealignedTemplate-After  $M = 0.8, SD = 0.11$ ;  $t(29) = 7.77, p < 0.001$ ;  $M_{diff} = 0.13, 95\% CI = [0.09, 0.16]$ ). Moreover, the differences between Template-Before and Template-After were more than four times larger in the Realigned heatmaps than in the Flipped ones (FlippedTemplate-After minus FlippedTemplate-Before:  $M = 0.03, SD = 0.10$ ; RealignedTemplate-After minus RealignedTemplate-Before:  $M =$

716 0.13, SD = 0.09), and the difference between these differences was statistically significant  
717 ( $t(29) = 3.81$ ,  $p < 0.001$ ;  $M_{diff} = 0.10$ , 95% CI = [0.05, 0.15]).

718         Similar to Experiment 1, we conducted an analysis of the proportion of fixations landing  
719 within flipped and re-aligned ROIs on the two-tone images to assess in more detail whether  
720 fixations are specifically object-oriented. Note, however, that to the extent to which ROIs cross  
721 the central vertical axis of an image, flipped ROIs overlap with re-aligned ROIs (this happened  
722 in 16 images, with an average overlap of 49.59 % (SD = 29.16 %) of pixels). To ensure that  
723 ROIs are unique, in this analysis, we used flipped and re-aligned ROIs from which the overlap  
724 between the two had been removed. The proportion of fixations landing in the flipped ROIs did  
725 not differ between the After and the Before conditions (Before:  $M = 7\%$ ,  $SD = 7$ ; After:  $M = 7\%$ ,  
726  $SD = 7$ ;  $t(29) = 0.14$ ,  $p = 0.888$ ;  $M_{diff} = 0$ , 95% CI = [-1, 2]). The same metric for the realigned  
727 ROIs indicated a clear difference between the two conditions, with more fixations landing in the  
728 realigned ROI in the After than the Before condition, indicating that changes in fixations were  
729 object-specific (Before:  $M = 16\%$ ,  $SD = 11$ ; After:  $M = 19\%$ ,  $SD = 11$ ;  $t(29) = 3.55$ ,  $p < 0.01$ ;  
730  $M_{diff} = 4\%$ , 95% CI = [2, 6]).

731         Similar to the findings for all fixations, heatmap similarities did not differ when comparing  
732 the FlippedTemplate-Before pair vs. the FlippedTemplate-After pair for first fixations  
733 (FlippedTemplate-Before:  $M = 0.44$ ,  $SD = 0.25$ ; FlippedTemplate-After:  $M = 0.45$ ,  $SD = 0.27$ ;  
734  $t(29) = 0.47$ ,  $p = 0.645$ ;  $M_{diff} = 0.01$ , 95% CI = [-0.03, 0.05]). By contrast to all fixations,  
735 however, the equivalent comparison for the realigned pairs did also not show a significant  
736 difference, albeit with a numerically larger effect in the direction expected from Experiment 1  
737 (RealignedTemplate-Before:  $M = 0.57$ ,  $SD = 0.21$ ; RealignedTemplate-After  $M = 0.62$ ,  $SD =$   
738  $0.20$ ;  $t(29) = 1.87$ ,  $p = 0.072$ ;  $M_{diff} = 0.05$ , 95% CI = [0, 0.09]).

739         The ROI analyses for first fixations corroborated this pattern of results. We found no  
740 significant differences between the After and the Before conditions in the proportion of fixations  
741 landing in the flipped ROI (Before:  $M = 7\%$ ,  $SD = 10$ ; After:  $M = 6\%$ ,  $SD = 8$ ; Before-After:  $t(29)$

742 = -0.87,  $p = 0.391$ ;  $M_{diff} = -1$ , 95% CI = [-4, 2]) and the realigned ROI (Before:  $M = 13\%$ ,  $SD =$   
743  $18$ ; After:  $M = 16\%$ ,  $SD = 17$ ; Before-After:  $t(29) = 1.66$ ,  $p = 0.107$ ;  $M_{diff} = 3$ , 95% CI = [-1, 6]),  
744 albeit with a numerical pattern in line with that of all fixations.

745

### 746 **Comparison between Experiments 1 and 2**

747 The spatial misalignment of template and two-tone images had an influence on how well  
748 observers were able to disambiguate the two-tones, as indicated by the finding that the (per-  
749 image) average increase in the meaningfulness ratings in Experiment 2 were smaller than in  
750 Experiment 1 ( $t(29) = 8.63$ ,  $p < 0.001$ ;  $M_{diff} = 0.12$ , 95% CI = [0.09, 0.15]). In order to contrast  
751 the effects on gaze guidance across experiments, we directly compared the increase in  
752 similarity between the Template-Before vs. Template-After pairs across Experiments 1 and 2.  
753 Given that both experiments differed with respect to the number of observers who contributed  
754 to the heatmaps of each image, we included fixations only from 18 observers from Experiment  
755 1 (drawn randomly). We found that the increase in similarity between the Template-Before vs.  
756 Template-After pairs were larger in Experiment 1 than Experiment 2 (Experiment 1:  $M = 0.17$ ,  
757  $SD = 0.13$ ; Experiment 2:  $M = 0.13$ ,  $SD = 0.09$ ;  $p = 0.0174$ ;  $M_{diff} = 0.05$ , 95% CI = [0.01,  
758 0.08]). To ensure that the outcome did not depend on the specific set of observers from  
759 Experiment 1, we repeated this analysis for 20 different, randomly drawn sets and obtained the  
760 same pattern of outcomes for 19 of them.

761

### 762 **Experiment 2 – Discussion**

763 In sum, despite the spatial misalignment of objects in template and two-tone images,  
764 fixations were strongly influenced by object locations in Experiment 2. There was no evidence to  
765 suggest that mapping objects to locations played a role in gaze guidance. It is noteworthy,  
766 however, that the spatial misalignment between template and two-tone images in Experiment 2  
767 had an attenuating effect on the influence of objects on eye movements compared to

768 Experiment 1. Interestingly, this attenuation in gaze guidance data was mirrored by an  
769 attenuation in the meaningfulness ratings, reflecting the ability of observers to use prior  
770 knowledge to organise two-tone images into meaningful object percepts (which was,  
771 nevertheless, robust). This finding is consistent with our overall interpretation that knowledge-  
772 driven object representations are important in eye-movement control.

773         While the analysis of first fixations showed a pattern that was numerically similar to that  
774 of all fixations, none of the analyses reached significance. In other words, by contrast to  
775 Experiment 1, first fixations in Experiment 2 did not show significant object-oriented effects,  
776 probably because the spatial misalignment between template and two-tone images resulted in  
777 less efficient perceptual organization of the latter into a meaningful percept (as suggested by the  
778 comparison of the meaningfulness ratings between Experiments 1 and 2). Importantly, analyses  
779 of first fixations also provided no evidence to suggest that a process of object-to-location  
780 mapping played any role in guiding first fixations during viewing of the two-tone images. Taken  
781 together, the results from Experiment 2 exclude the possibility that gaze guidance in the After  
782 condition is based on a mapping of objects to locations via retrieval of this information from the  
783 Template condition.

784         In a third experiment, we addressed two further alternative explanations of the results  
785 from Experiment 1. First, it is possible that during the phase when two-tone images are blended  
786 with templates, observers learn to associate specific image-features in the two-tone images with  
787 object locations in the templates. When viewing two-tone images in the After condition, these  
788 feature-object associations might guide fixations towards these specific visual patterns,  
789 irrespective of transformations such as those introduced by the mirror-flipping. While this  
790 possibility might seem implausible, there is evidence to suggest that such learning processes  
791 are an important factor in oculomotor control (Alfandari, Belopolsky, & Olivers, 2019).

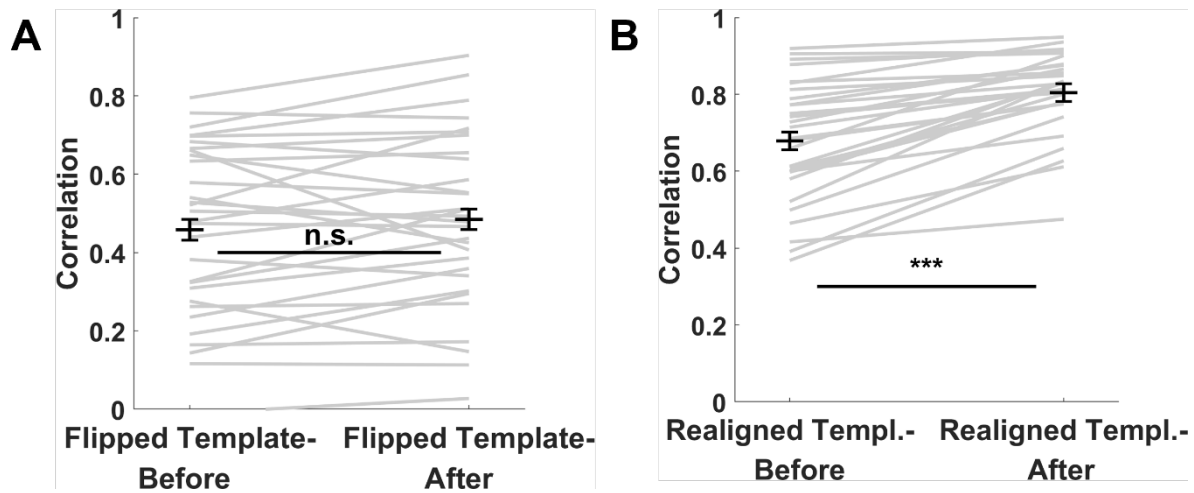
792         A final alternative explanation of our results from both Experiment 1 and 2 relates to  
793 potential order effects. It is possible that the changes in fixation patterns between Before and

794 After conditions resulted from viewing two-tones for a second time, rather than from knowledge-  
 795 dependent perceptual organization. In other words, observers might sample information from  
 796 different image regions on second compared to first viewing, irrespective of the kind of  
 797 information they acquire in the meantime. We conducted Experiment 3 to exclude the possibility  
 798 that (i) feature-object associations, or (ii) any order effects could explain the effects of  
 799 Experiments 1 and 2.

800

801 **Figure 10**

802 *Results of Experiment 2*



803

804 *Note.* A) Similarities between heatmaps from two-tone images and mirror-flipped templates,  
 805 where the two-tones were viewed either in the Before or in the After condition. The heatmaps  
 806 derived from the mirror-flipped template images were used either before (A) or after (B) the  
 807 mirror-flipping was reverted by ‘flipping back’ these heatmaps and realigning them with the  
 808 heatmaps from two-tone images.

809

810 **Experiment 3**

811 **Overview**

812 Experiment 3 adopted the same procedure as the previous experiments except that the

813 templates from Experiment 1 ('real templates') were replaced with different images that were  
814 unrelated to the two-tones ('dummy templates'). This experimental design allowed us to test  
815 whether feature-object associations provide a plausible explanation for the findings of  
816 Experiment 1 and 2. Specifically, observers might associate certain features in the two-tone  
817 images with objects in the templates during the Blending Phase. When viewing two-tone images  
818 in the After condition, these feature-object associations could drive fixations towards image  
819 locations in the two-tones that overlap with objects in the respective (dummy) templates. These  
820 effects should be observable despite observers not having acquired the prior object-knowledge  
821 required to organize the two-tone images into coherent percepts. Moreover, the design also  
822 allowed us to assess whether order effects could explain the findings from Experiments 1 and 2.

823

### 824 **Experiment 3 – Method**

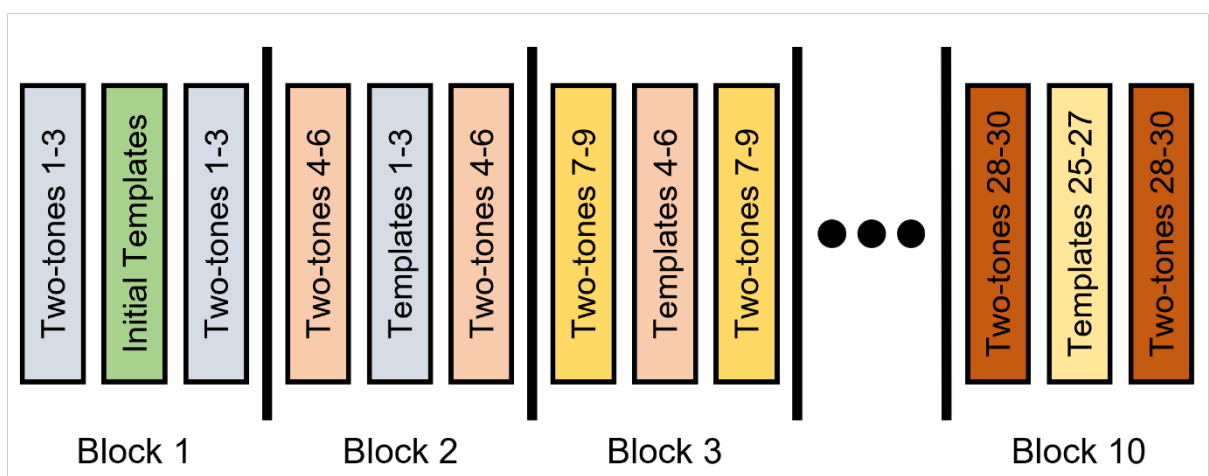
825 Experiment 3 was completed by 20 observers (mean age 19.55, 5 men) who did not  
826 participate in the previous two experiments. All were Cardiff University students. The  
827 procedure was identical to the previous experiments except that in each block, the templates  
828 used in the Template condition and in the Blending Phase were unrelated to the two-tones  
829 presented in this block ('dummy templates'). Each two-tone had a unique dummy template  
830 paired with it and this pairing was fixed for all observers. Importantly, each dummy template  
831 was a 'real template' of a different two-tone presented in the preceding block during the  
832 experiment (see Fig. 11). While templates in this experiment could thus not provide object  
833 knowledge that would help organize the two-tone image into an object percept in the After  
834 condition, we were nevertheless able to register eye movements on the real templates.  
835 Measuring fixations on real templates was necessary to assess whether simply viewing a two-  
836 tone for a second time, without prior object-knowledge, would lead to increased similarity  
837 between heatmaps of two-tone images in the After and their real templates, as seen in the  
838 previous experiments.

839 In the first block, the same dummy templates – greyscale images not related to any of  
 840 the two-tones – were always presented. In all other blocks, the assignment of stimuli to  
 841 experimental blocks was pseudo-randomized for each observer individually in a way which  
 842 guaranteed that dummy templates presented in any given block were the real templates of two-  
 843 tones presented in the preceding block (see Fig. 11). To ensure that we included data from the  
 844 same number of observers for each two-tone and template, we had to discard fixations  
 845 registered on the two-tones presented in the final experimental block and fixations from the  
 846 dummy templates from the first blocks ('initial templates'). Note that – because we pseudo-  
 847 randomized the order of stimulus presentation for each observer individually – for different  
 848 images we had to discard data from different observers. Importantly, however, for each image  
 849 set consisting of a two-tone (viewed in Before and After condition), its dummy template, and its  
 850 real template, we retained data from a homogenous group of 18 observers (out of 20 who  
 851 completed the experiment), but the composition of these groups was different for different image  
 852 sets.

853

854 **Figure 11**

855 *Randomization schema used in Experiment 3*



856

857 *Note.* Within each block, stimuli were presented in a randomised order (as in Experiment 1 and

858 2). The presentation of images was arranged in such a way that templates in, e.g., Block 2,  
859 were the real templates of the two-tone images in Block 1. This order allowed us to register  
860 fixations for the real templates (for comparison with fixation on two-tone images) while omitting  
861 the opportunity for the observer to acquire the relevant prior object-knowledge that would allow  
862 them to disambiguate the two-tone images.

863

864

### Experiment 3 – Results

865

#### Analysis of meaningfulness ratings

866

867

868

869

870

871

872

873

874

875

#### Controlling for the effects of object-to-feature mapping

876

877

878

879

880

881

882

883

Experiment 3 tested the hypothesis that the effects observed in the two previous experiments might be explainable by a learned association between feature clusters in two-tones and object locations on templates. Specifically, it is possible that during blending of two-tone images and templates, observers learn to associate specific features of the two-tones with object locations in the templates and then re-visit these features when viewing the two-tone images in the After condition. Our analysis indicated that the similarity in heatmaps in the DummyTemplate-After pair was higher compared to the DummyTemplate-Before pair (Fig. 12C). This increase in similarity, although significant in a statistical sense, was small



884 (DummyTemplate-Before:  $M = 0.46$ ,  $SD = 0.21$ ; DummyTemplate-After:  $M = 0.52$ ,  $SD = 0.22$ ;  
885  $t(29) = 4.70$ ,  $p < 0.001$ ;  $M_{diff} = 0.06$ ,  $95\% CI = [0.03, 0.08]$ ). Nevertheless, the analysis  
886 provided evidence to suggest that feature-object associations might guide oculomotor control  
887 to a limited extent. Alternatively, these results could be driven by memory retrieval of object-  
888 locations in the templates: while Experiment 2 showed that memory retrieval does not play a  
889 role when perceptual organization takes place, this process may become important when the  
890 stimulus remains unorganized with no object representations to guide eye movements. In  
891 either case, it is interesting that the analysis of first fixations did not indicate a difference  
892 between DummyTemplate-Before and DummyTemplate-After (DummyTemplate-Before:  $M =$   
893  $0.41$ ,  $SD = 0.21$ ; DummyTemplate-After:  $M = 0.45$ ,  $SD = 0.24$ ;  $t(29) = 1.38$ ,  $p = 0.179$ ;  $M_{diff} =$   
894  $0.04$ ,  $95\% CI = [-0.02, 0.01]$ ). This finding suggests a different temporal development of the  
895 influence on gaze guidance by object representations vs. by object-to-location or object-to-  
896 feature mappings: while the former is present from the first fixations, the latter kick in only after  
897 the first fixation (and potentially only if no object representations are available to provide  
898 guidance).

899 The ROI analyses corroborated the findings for heatmaps: for all fixations, we found a  
900 significant difference between the After and the Before conditions in the proportion of fixations  
901 landing in the ROIs of DummyTemplates (Before:  $M = 0.21$ ,  $SD = 0.24$ ; After:  $M = 0.23$ ,  $SD =$   
902  $0.27$ ; Before-After:  $t(29) = 2.64$ ,  $p = 0.013$ ;  $M_{diff} = 0.02$ ,  $CI = [0.01, 0.04]$ ). Note that this  
903 difference was similar in magnitude to the equivalent difference regarding the ROIs of real  
904 Templates (see the *Controlling for order effects* section; difference between the differences:  
905  $t(29) = -1.27$ ,  $p = 0.212$ ;  $M_{diff} = -0.02$ ,  $95\% CI = [-0.05, 0.01]$ ). Finally, first fixations showed no  
906 difference in the proportion of fixations landing in the ROIs of the DummyTemplates (Before:  $M$   
907  $= 0.26$ ,  $SD = 0.34$ ; After:  $M = 0.27$ ,  $SD = 0.35$ ; Before-After:  $t(29) = 0.79$ ,  $p = 0.435$ ;  $M_{diff} =$   
908  $0.01$ ,  $95\% CI = [-0.02, 0.05]$ ).

909

**910 Object-to-location mapping: comparison between Experiments 1 and 3**

911 While the results reported in the previous section suggest that object-to-location or  
912 object-to-feature mapping might influence gaze guidance in the After condition (after the first  
913 fixation), the key question is whether these effects can explain the results found in Experiment  
914 1. To address this issue, we directly compared the increase in similarity between the Template-  
915 Before vs. Template-After pairs across Experiments 1 and 3. Given that both experiments  
916 differed with respect to the number of observers who contributed to the heatmaps of each  
917 image, we adopted a similar approach for that used to compare Experiments 1 and 2 (i.e., we  
918 randomly drew 18 observers from Experiment 1 and repeated this analysis for 20 different,  
919 randomly drawn sets). This analysis indicates that the change in similarity between the  
920 Template-Before vs. Template-After pairs was larger in Experiment 1 than in Experiment 3  
921 (Experiment 1:  $M = 0.17$ ,  $SD = 0.13$ ; Experiment 3:  $M = 0.06$ ,  $SD = 0.07$ ;  $t(29) = 4.15$ ,  $p < 0.001$ ;  
922  $M_{diff} = 0.11$ , 95% CI = [0.06, 0.17]; results for one of the 20 sets).

923 Our results (for all fixations) thus demonstrate that the processes responsible for  
924 changing gaze-patterns between Before and After conditions in Experiment 3 cannot fully  
925 explain the analogous changes in Experiment 1. One possible explanation for this finding is that  
926 it might be more difficult to learn object-to-feature mappings in Experiment 3 than Experiment 1  
927 (during viewing of the template images and the blending phase). If we assume that gaze is  
928 guided by this mapping process, then less robust learning might explain the differences in effect  
929 size for all fixations in Experiments 1 and 3. Importantly, however, the differences in temporal  
930 trajectories found in the two experiments might be difficult to reconcile with this idea: by contrast  
931 to Experiment 1, we found no evidence for a change between Before and After in Experiment 3  
932 for first fixations. This pattern of results suggests that (partly) different processes that are  
933 characterised by different temporal trajectories are at work in the two experiments. Specifically,  
934 we argue that the influence of object representations is present from the first fixations onwards  
935 (as seen in Experiment 1), while object-to-feature or object-to-location mapping kicks in later (as

936 seen in Experiment 3), and potentially only if no object representations are available to provide  
937 guidance. Overall, the pattern of results in Experiments 1 and 3 suggest that the findings for first  
938 fixations cannot be explained by either object-to-feature or object-to-location mapping, even if  
939 these processes might contribute to, but not fully explain, the effect seen in all fixations.

940

### 941 **Controlling for order effects**

942 In the final analysis, we considered the possibility that order effects explain the key  
943 findings of Experiment 1 and 2. Specifically, we asked whether viewing the same two-tones for  
944 a second time without receiving prior object-knowledge could change fixation patterns such that  
945 they would resemble the patterns from the (real) templates. Recall that the design of Experiment  
946 3 ensured that observers saw each two-tone image twice, each time without prior object-  
947 knowledge (Before and After conditions, respectively) and they also saw the real template for  
948 these two-tones in the following block. If the findings in Experiments 1 and 2 resulted, at least  
949 partly, from an order effect, we would expect that the similarity in fixation patterns in the (real)  
950 Template-After pair would be higher than in the (real) Template-Before pair in the current  
951 experiment.

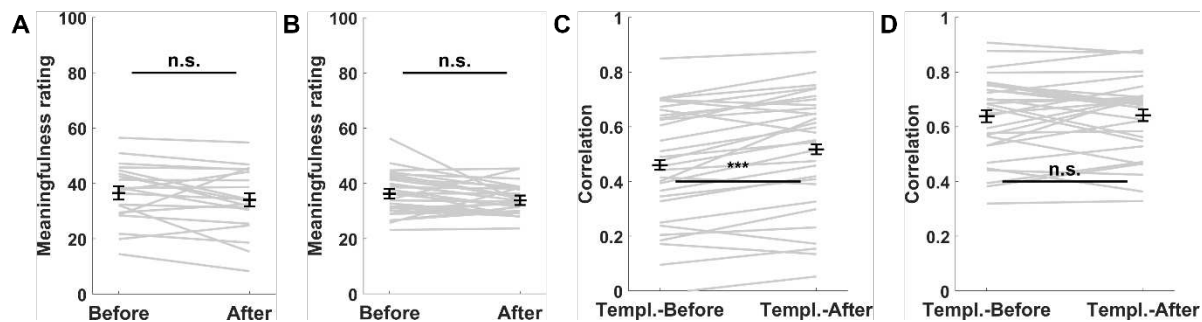
952 The results were inconsistent with this 'second-viewing' hypothesis (Fig. 12D). The  
953 heatmap similarities between the real templates and the corresponding two-tones viewed in the  
954 Before and After conditions were not statistically different (Template-Before  $M = 0.64$ ,  $SD =$   
955  $0.15$ ; Template-After  $M = 0.64$ ,  $SD = 0.14$ ;  $t(29) = 0.22$ ,  $p = 0.830$ ;  $M_{diff} = 0$ , 95% CI = [-0.03,  
956 0.03]). Moreover, a Bayes factor analysis provided evidence to support a lack of a difference  
957 (BF = 0.20). We found a similar result for first fixations (Template-Before  $M = 0.52$ ,  $SD = 0.21$ ;  
958 Template-After  $M = 0.53$ ,  $SD = 0.18$ ;  $t(29) = 1.01$ ,  $p = 0.323$ ;  $M_{diff} = 0.02$ , 95% CI = [-0.02,  
959 0.06]). Finally, the ROI analyses for both all fixations (Before:  $M = 0.27$ ,  $SD = 0.24$ ; After:  $M =$   
960  $0.27$ ,  $SD = 0.33$ ; Before-After:  $t(29) = 0.32$ ,  $p = 0.212$ ;  $M_{diff} = -0.02$ , 95% CI = [-0.05, 0.01]) and  
961 first fixations corroborated these findings (Before:  $M = 0.31$ ,  $SD = 0.34$ ; After:  $M = 0.31$ ,  $SD =$

962 0.33; Before-After:  $t(29) = 0.44$ ,  $p = 0.666$ ;  $M_{diff} = 0.01$ , 95% CI = [-0.03, 0.05]).

963

964 **Figure 12**

965 *Results of Experiment 3*



966

967 *Note.* Meaningfulness ratings averaged per observer (A) and per image (B). C) Comparison of  
 968 heatmap similarities between two-tones (viewed in the Before and After conditions) and their  
 969 dummy templates (i.e., unrelated images). D) Comparison of heatmap similarities between two-  
 970 tones (viewed in Before and After conditions) and their real templates.

971

972

### Discussion

973

974

975

976

977

978

979

980

981

982

When an observer explores the environment with no specific task other than to obtain information, eye movements are typically directed towards object locations. Here, we consider this effect in light of emerging evidence highlighting the complex and intricate relationship between image-computable features and high-level object representations in visual perception. Specifically, we ask whether object-oriented eye-movements result from gaze being guided by high-level features or by objecthood, i.e., the fact that these features are bound into an object representation. We recorded eye movements in response to two-tone images, stimuli that appear as meaningless patches on initial viewing but, once relevant object-knowledge has been acquired, are organized into coherent and meaningful percepts of objects. In the current study, prior object-knowledge was provided in the form of template images, i.e., the unambiguous

983 photographs from which the two-tone images had been generated. Across three experiments,  
984 fixation patterns on the same two-tone images differed substantially depending on whether  
985 observers experienced them as meaningless patches or organized them into object  
986 representations. In particular, when organized into object representations, we found that fixation  
987 patterns on two-tone images were more similar to those on templates, more focused on object-  
988 specific, pre-defined regions-of-interest, less dispersed, and more consistent across observers.  
989 These effects were evident from the first fixations on an image. Importantly, eye-movements on  
990 two-tone images were best explained by a simple model that takes into account both low-level  
991 features and high-level, knowledge-dependent object representations. Together, these findings  
992 highlight the importance of dynamic interactions between image-computable features and  
993 knowledge-driven perceptual organization in guiding information sampling via eye-movements  
994 in humans.

995         The idea that knowledge-driven object representations restructure human eye-  
996 movements is supported by both our general assessment of fixation distributions between two-  
997 tone images and template, and also by a more specific analysis focusing on fixations within  
998 regions-of-interest. These findings provide strong support for the hypothesis that objecthood per  
999 se contributes to the process of selecting fixation targets in images. In our experimental design,  
1000 image-computable visual features are insufficient for object representations to emerge, their  
1001 formation is dependent on prior object-knowledge. This characteristic of two-tone images is an  
1002 important experimental tool: it allows us to decisively rule out the possibility that human  
1003 oculomotor control during free viewing relies solely on image-computable features, regardless of  
1004 whether these features are low- or high-level (Zelinsky & Bisley, 2015). The simple but critical  
1005 result in this regard is the finding that eye-movement patterns differed dependent on whether  
1006 observers had formed object representations despite the fact that the features in the stimuli  
1007 remained identical. Of course, despite being highly impoverished, two-tone images might still  
1008 contain some of the features that give rise to object representations in the Template images.

1009 Note, however, that Before and After conditions have identical featural overlap with the  
1010 Template condition, and differences in eye-movements between Before and After can therefore  
1011 not be explained by this factor.

1012 In addition to its use as an experimental tool, however, the dependence of object  
1013 representations on prior knowledge is also important from a conceptual perspective.  
1014 Specifically, the finding that fixations were guided by knowledge-dependent representations  
1015 demonstrates that for the oculomotor system, objects cannot be conceptualised (exclusively) as  
1016 image-computable, high-level features (Schütt et al., 2019). As highlighted in the introduction,  
1017 Schütt and colleagues' (2019) study is one of the few that is explicit about this  
1018 conceptualisation. While other studies have been less clear about exactly what constitutes an  
1019 object, many treat them in a manner that (implicitly) equates object representations to complex  
1020 high-level features (Borji & Tanner, 2016; Einhäuser et al., 2008; Nuthmann et al., 2020; Pajak  
1021 & Nuthmann, 2013; Stoll et al., 2015). While these studies contribute to our understanding of the  
1022 role of low- vs. high-level features in gaze control, they are not able to (and did not intend to)  
1023 dissociate the influence of image-computable features from the influence of objecthood per se.  
1024 Here, we show that objecthood that is relevant for guiding eye-movements is a characteristic  
1025 that is distinct from the collection of any low- or high-level features. In our study, objecthood  
1026 emerges in the interaction between prior object-knowledge and the visual input. Whether object  
1027 representations that are relevant for oculomotor control are always distinct from the featural  
1028 input is a difficult question that we cannot answer with our data. However, the size, the speed,  
1029 and the incidental nature of these effects suggests that they might be characteristic of eye-  
1030 movement control in everyday visual behaviour.

1031 Our findings contrast in interesting ways with previous work that studied the relationship  
1032 between eye-movements and object representations using ambiguous, bi-stable object stimuli  
1033 (Kietzmann et al. 2011, 2015). These studies demonstrate that fixation patterns typical for one  
1034 of the two interpretations of these stimuli often precede the emergence of the first percept

1035 corresponding to that interpretation. Thus, eye movements might play a role in the accumulation  
1036 of image-computable evidence for competing stimulus interpretations, potentially suggesting  
1037 that specific fixation patterns facilitate selection of one of two possible interpretations. In  
1038 contrast to this finding, our results suggest that the influence of object representations precedes  
1039 the first saccade. While our data provide no means to reconcile these contrasting findings, one  
1040 possibility is a bi-directional relationship, where object representations guide eye-movements  
1041 (as shown here) and eye-movements also support the generation of object representations (as  
1042 shown in the studies by Kietzmann and colleagues). The use of a design that focusses on the  
1043 role of eye movements in the accumulation of image-computable evidence for competing  
1044 stimulus interpretations might be the reason why Kietzmann and colleagues mainly picked up  
1045 on the latter component.

1046         Manipulating low-level features is another approach aiming at dissociating feature-based  
1047 and object-based effects. It was adopted by Stoll and colleagues (2015), who reduced contrast  
1048 – a low-level feature contributing to saliency – in image areas containing objects. Given that in  
1049 this study, objects are defined by high-level features, this approach provides a useful tool to  
1050 assess the influence of low- vs. high-level features. It does not, however, allow for distinguishing  
1051 between high-level features and objecthood per se as we do in the current study.

1052         Equally important as the finding that knowledge-driven object representations guide  
1053 human gaze is the fact that they do not fully determine the selection of fixation locations. While  
1054 eye-movements on two-tone images changed once they elicited object representations such  
1055 that fixation distributions became more similar to fixations on template images, substantial  
1056 differences in eye-movements remained between these two conditions. Our linear combination  
1057 analysis suggests that this disparity is systematic and can be explained by the differences in the  
1058 features in two-tone vs. template images. In this analysis, we generated linear combinations  
1059 with varying proportions of the heatmaps from the Template and Before conditions. We then  
1060 assessed the similarities between these combined heatmaps and the heatmaps from the After

1061 condition. These similarities peaked for combined heatmaps that were determined by the  
1062 fixation distributions from both the Template and the Before conditions (and not just one of  
1063 them). The finding thus demonstrates that when observers experienced the percept of an object  
1064 in the two-tone images (After condition), fixations were best explained by a combination of the  
1065 factors guiding eye movements in the Before and the Template conditions. Specifically, even  
1066 when observers perceived an object in the two-tone images, their eye movements were only  
1067 partly determined by the factors that guide eye movements in response to the template image.  
1068 The image-computable features that drive eye movements in response to two-tone images  
1069 when no object is perceived (Before condition) still made a substantial contribution to gaze  
1070 guidance. Note that the linear combination analysis was conducted on a per-image basis. The  
1071 finding that both features and objecthood contribute to eye-movement control can therefore not  
1072 be explained by averaging across different images, with some leading to purely feature-driven  
1073 and other to purely representation-driven eye-movement control.

1074         The finding that features remain important for eye-movement control even after having  
1075 been bound into a high-level object representation potentially challenges some of the strong  
1076 claims regarding the role of features vs. objects in gaze guidance. For instance, the cognitive  
1077 relevance theory (Henderson et al., 2009) proposes that visual features do not contribute to  
1078 oculomotor control directly but provide the means to generate a representation of potential  
1079 fixation locations that have not yet been ranked for priority. High-level factors operate on this  
1080 'flat landscape' to determine the ultimate fixation locations. In other words, features are  
1081 important only as potential carriers of higher-level representations and do not contribute to eye-  
1082 movement control by themselves. According to this idea, as long as visual features give rise to  
1083 similar object representation, these representations should guide eye movements towards  
1084 similar locations, independently of the specific characteristics of features. Therefore, to the  
1085 extent to which two-tones and templates lead to similar object representations, both image  
1086 types should result in similar eye-movement patterns independent of their featural differences.



1087 Contrasting with this notion, in the analysis of linear combinations, we found that the specific  
1088 features that support these high-level representations continue to exert a sizeable influence on  
1089 eye-movements. Specifically, we demonstrate that the same features that guided eye-  
1090 movements when no object representation was present (Before condition) still had an influence  
1091 on gaze guidance when an object representation had been generated (After). Therefore, to the  
1092 extent to which two-tones and templates lead to similar object representations, we would have  
1093 expected both image types to result in similar eye-movement patterns independent of their  
1094 featural differences. Contrasting with this notion, we found that, while features can be flexible  
1095 carriers of object representations that guide eye-movements as predicted by the cognitive  
1096 relevance theory, the specific features that support these high-level representations persist to  
1097 exert a sizeable influence.

1098         In terms of the time-course of eye-movements, we provide clear evidence that already  
1099 the first fixations after image onset are affected by objecthood. Interestingly, however, the linear  
1100 combination analysis indicates that for first fixations the relative influence of features is stronger  
1101 – and, therefore, the relative influence of objecthood weaker – compared to later fixations. Thus,  
1102 while the influence of knowledge-dependent object representations emerges quickly, the linear  
1103 combination analysis suggests that the effects of knowledge-driven perceptual organization  
1104 continue to build beyond the first fixation, by contrast to the effects of features. Nevertheless,  
1105 our data suggest that the influence of knowledge-dependent object representations emerges  
1106 quickly and exerts an influence from the earliest fixations.

1107         At image onset, when the eyes are stationary prior to the first saccade, most of the  
1108 image is viewed via peripheral vision with only a small part being inspected with high-resolution  
1109 foveal vision. The analysis of the first fixations therefore suggest that the visual system is able to  
1110 generate knowledge-dependent object representations quickly and largely based on information  
1111 from peripheral vision. Due to the optical, anatomical, and neurophysiological characteristics of  
1112 the primate visual system, peripheral vision is limited in various respects (Rosenholtz, 2016),

1113 but there is good evidence that it provides enough information to generate a gist representation  
1114 of a visual scene that can guide subsequent eye movements (Anderson, Donk, & Meeter, 2016;  
1115 Castelhana & Henderson, 2007; Melissa L.H. Võ & Schneider, 2010). Exactly how detailed this  
1116 gist representation is, which features it contains, and whether objects are represented varies  
1117 depending on a number of different factors (Malcolm, Groen, & Baker, 2016; Wallis, Bethge, &  
1118 Wichmann, 2016). Note, however, that this question is of limited relevance in the current context  
1119 because features in two-tone images – independently of whether they are viewed by foveal or  
1120 peripheral vision – are necessary but, by themselves, not sufficient to determine the high-level  
1121 object representations we study here. However, one notion that might help in explaining the  
1122 rapid influence of knowledge-dependent object representations on eye movements is provided  
1123 by the suggestion that object recognition involves a predictive process that is triggered by low  
1124 spatial-frequencies in the input (Bar et al., 2006; Bar, 2003, 2004, 2021; Bullier 2001).  
1125 Specifically, low spatial-frequency information is thought to be fed forward by fast projections to  
1126 high-level brain systems that connects this rudimentary input to prior object-knowledge. This  
1127 process narrows down the search space of possible hypotheses about object identities in the  
1128 input, thereby scaffolding and shaping a more precise perceptual experience of the input. It is  
1129 therefore tempting to speculate that, in our experiment, first fixations were guided by object  
1130 representations that are based on the process that links impoverished low spatial-frequency  
1131 image content to prior knowledge, while later fixations might be based on fuller object  
1132 representations. This idea rests on the assumption that two-tone images provide low spatial-  
1133 frequency information to peripheral vision that allows the linking of two-tone images to memory  
1134 representations of template images. Given that the image-processing operations required to  
1135 generate two-tone images mainly affect high spatial-frequency components and have less  
1136 impact on low spatial frequencies, this assumption seems plausible.

1137           While our analyses mainly focused on locations of fixations, other aspects of oculomotor  
1138 control are also influenced by knowledge-dependent perceptual organization. Specifically, we

1139 observed a decrease in saccade length and an increase in fixation duration when two-tone  
1140 images were organized into object representations (After condition) compared to when they  
1141 were not (Before condition). Both changes are indicative of a shift from image exploration to  
1142 image exploitation (Gameiro et al., 2017; Kaspar et al., 2013), an interpretation that was also  
1143 supported by the decrease in entropy across the two conditions. The oculomotor system  
1144 constantly has to decide whether to keep the eyes still in order to be able to further inspect the  
1145 currently fixated scene region – a process referred to as exploitation –, or to perform a saccade  
1146 to explore another part of the image. Interestingly, in our study, the shift from exploration to  
1147 exploitation went along with an increase in the amount of fixations landing on objects. This  
1148 finding suggests that the visual system prioritizes objects in a specific way: it exploits object  
1149 locations for further information while abandoning exploration of the remaining parts of the  
1150 image. In other words, our data demonstrate that clusters of features that are bound into, and  
1151 provide support for, object representations become interesting for the visual system over non-  
1152 object related feature clusters (for a similar finding, see Król & Król, 2019). The shift from  
1153 exploitation to exploration once objecthood is established also leads to higher consistency  
1154 across observers. This finding suggests that guidance of exploration is either more idiosyncratic  
1155 or that image-computable features that are not bound into object representations do not provide  
1156 strong constraints for oculomotor control. Conversely, object representations, even when  
1157 supported by exactly the same features, have a structuring or normative effect on information  
1158 sampling. In other words, while observers explore features in different ways, they exploit objects  
1159 in similar ways.

1160           In summary, we demonstrate that gaze guidance is best understood by dynamic  
1161 interactions between image-computable features and knowledge-dependent perceptual  
1162 organization. Specifically, our findings demonstrate the importance of objecthood per se – i.e.,  
1163 representations that are not reducible to image-computable features – in oculomotor control but  
1164 also indicate a persistent contribution of object-independent features. We demonstrate that

1165 when visual input remains identical, the emergence of knowledge-dependent object  
1166 representations substantially restructures information sampling via eye-movements. However,  
1167 we also show that even when image-computable features are bound into object representations,  
1168 they still retain some influence on eye movements, challenging the idea that the role of features  
1169 is limited to being carriers for high-level representation without direct influence on eye-  
1170 movements. Finally, we also show that the emergence of object representations results in an  
1171 overall change of the information-sampling strategy of the visual system, leading to the  
1172 prioritization of information extraction from features that are bound into object representations,  
1173 at the expense of exploration of the entire image.

1174

1175 **CRedit authorship statement:**

1176 M.P.: Conceptualisation, Methodology, Software, Validation, Formal analysis, Investigation,

1177 Visualization, Writing - Original Draft, Writing - Review & Editing

1178 E. v.d. H.: Methodology, Writing - Review & Editing

1179 C.T.: Conceptualisation, Methodology, Writing - Original Draft, Writing - Review & Editing,

1180 Supervision, Resources

**References**

- 1181  
1182 Alfandari, D., Belopolsky, A. V., & Olivers, C. N. L. (2019). Eye movements reveal learning and  
1183 information-seeking in attentional template acquisition. *Visual Cognition*, 27(5–8), 467–486.  
1184 <https://doi.org/10.1080/13506285.2019.1636918>
- 1185 Anderson, N. C., Donk, M., & Meeter, M. (2016). The influence of a scene preview on eye  
1186 movement behavior in natural scenes. *Psychonomic Bulletin & Review*, 1–8.  
1187 <https://doi.org/10.3758/s13423-016-1035-4>
- 1188 Anderson, N. C., Ort, E., Kruijine, W., Meeter, M., & Donk, M. (2015). It depends on when you  
1189 look at it: Saliency influences eye movements in natural scene viewing and search early in  
1190 time. *Journal of Vision*, 15(5), 9. <https://doi.org/10.1167/15.5.9>
- 1191 Bar, M. (2003). A cortical mechanism for triggering top-down facilitation in visual object  
1192 recognition. *Journal of Cognitive Neuroscience*, 15(4), 600–609.  
1193 <https://doi.org/10.1162/089892903321662976>
- 1194 Bar, M. (2004). Visual objects in context. *Nature Reviews Neuroscience*, 5(8), 617–629.  
1195 <https://doi.org/10.1038/nrn1476>
- 1196 Bar, M., Kassam, K. S., Ghuman, A. S., Boshyan, J., Schmid, A. M., Dale, A. M., ... Halgren, E.  
1197 (2006). Top-down facilitation of visual recognition. *Proceedings of the National Academy of*  
1198 *Sciences*, 103(2), 449–454. <https://doi.org/10.1073/pnas.0507062103>
- 1199 Bar, M. (2021). From Objects to Unified Minds. *Current Directions in Psychological Science*,  
1200 30(2), 129–137. <https://doi.org/10.1177/0963721420984403>
- 1201 Borji, A., Sihite, D. N., & Itti, L. (2013). Objects do not predict fixations better than early saliency:  
1202 A re-analysis of einhäuser et al.'s data. *Journal of Vision*, 13(10), 1–4.  
1203 <https://doi.org/10.1167/13.10.18>
- 1204 Borji, A., & Tanner, J. (2016). Reconciling Saliency and Object Center-Bias Hypotheses in  
1205 Explaining Free-Viewing Fixations. *IEEE Transactions on Neural Networks and Learning*  
1206 *Systems*, 27(6), 1214–1226. <https://doi.org/10.1109/TNNLS.2015.2480683>

- 1207 Brainard, D. (1997). The Psychophysics Toolbox. *Spatial Vision*, 10(4).
- 1208 Bullier, J. (2001). Integrated model of visual processing. *Brain Research Reviews*, 36(2-3), 96–  
1209 107. [http://doi.org/10.1016/S0165-0173\(01\)00085-6](http://doi.org/10.1016/S0165-0173(01)00085-6)
- 1210 Bylinskii, Z., Judd, T., Oliva, A., Torralba, A., & Durand, F. (2016). What do different evaluation  
1211 metrics tell us about saliency models? *IEEE Transactions on Pattern Analysis and Machine*  
1212 *Intelligence*, 41(3), 740–757. <https://doi.org/10.1109/TPAMI.2018.2815601>
- 1213 Castelhana, M. S., & Henderson, J. M. (2007). Initial Scene Representations Facilitate Eye  
1214 Movement Guidance in Visual Search. *Journal of Experimental Psychology: Human*  
1215 *Perception and Performance*, 33(4), 753–763. <https://doi.org/10.1037/0096-1523.33.4.753>
- 1216 Cerf, M., Paxon Frady, E., & Koch, C. (2009). Faces and text attract gaze independent of the  
1217 task: Experimental data and computer model. *Journal of Vision*, 9(12), 1–15.  
1218 <https://doi.org/10.1167/9.12.1>
- 1219 Christensen, J. H., Bex, P. J., & Fiser, J. (2015). Prior implicit knowledge shapes human  
1220 threshold for orientation noise. *Journal of Vision*, 15(9), 1–15.  
1221 <https://doi.org/10.1167/15.9.24>
- 1222 Clarke, A. D. F., & Tatler, B. W. (2014). Deriving an appropriate baseline for describing fixation  
1223 behaviour. *Vision Research*, 102, 41–51. <https://doi.org/10.1016/j.visres.2014.06.016>
- 1224 Cousineau, D. (2005). Confidence intervals in within-subject designs: A simpler solution to  
1225 Loftus and Masson's method. *Tutorials in Quantitative Methods for Psychology*, 1(1), 42–  
1226 45. <https://doi.org/10.20982/tqmp.01.1.p042>
- 1227 DiCarlo, J. J., Zoccolan, D., & Rust, N. C. (2012). How does the brain solve visual object  
1228 recognition? *Neuron*, 73(3), 415–434. <https://doi.org/10.1016/j.neuron.2012.01.010>
- 1229 Drewes, J., Trommershäuser, J., & Gegenfurtner, K. R. (2011). Parallel visual search and rapid  
1230 animal detection in natural scenes. *Journal of Vision*, 11(2), 1–21.  
1231 <https://doi.org/10.1167/11.2.20>
- 1232

- 1233 Einhäuser, W. (2013). Objects and saliency: Reply to Borji et al. *Journal of Vision*, 13(10), 20–  
1234 20. <https://doi.org/10.1167/13.10.20>
- 1235 Einhäuser, W., Spain, M., & Perona, P. (2008). Objects predict fixations better than early  
1236 saliency. *Journal of Vision*, 8(14), 18.1-26. <https://doi.org/10.1167/8.14.18>
- 1237 Elazary, L., & Itti, L. (2008). Interesting objects are visually salient. *Journal of Vision*, 8(3).  
1238 <https://doi.org/10.1167/8.3.3>
- 1239 Federico, G., & Brandimonte, M. A. (2019). Tool and object affordances: An ecological eye-  
1240 tracking study. *Brain and Cognition*, 135(May), 103582.  
1241 <https://doi.org/10.1016/j.bandc.2019.103582>
- 1242 Feldman, J. (2003). What is a visual object? *Trends in Cognitive Sciences*, 7(6), 252–256.  
1243 [https://doi.org/10.1016/S1364-6613\(03\)00111-6](https://doi.org/10.1016/S1364-6613(03)00111-6)
- 1244 Flounders, M. W., González-García, C., Hardstone, R., & He, B. J. (2019). Neural dynamics of  
1245 visual ambiguity resolution by perceptual prior. *ELife*, 8, 1–25.  
1246 <https://doi.org/10.7554/eLife.41861>
- 1247 Foulsham, T., & Kingstone, A. (2013). Fixation-dependent memory for natural scenes: An  
1248 experimental test of scanpath theory. *Journal of Experimental Psychology: General*,  
1249 142(1), 41–56. <https://doi.org/10.1037/a0028227>
- 1250 Gameiro, R. R., Kaspar, K., König, S. U., Nordholt, S., & König, P. (2017). Exploration and  
1251 Exploitation in Natural Viewing Behavior. *Scientific Reports*, 7(1), 1–23.  
1252 <https://doi.org/10.1038/s41598-017-02526-1>
- 1253 Gilbert, C. D., & Li, W. (2013). Top-down influences on visual processing. *Nature Reviews*  
1254 *Neuroscience*, 14(5), 350–363. <https://doi.org/10.1038/nrn3476>
- 1255 Groen, I. I. A., Ghebreab, S., Prins, H., Lamme, V. A. F., & Scholte, H. S. (2013). From Image  
1256 Statistics to Scene Gist: Evoked Neural Activity Reveals Transition from Low-Level Natural  
1257 Image Structure to Scene Category. *Journal of Neuroscience*, 33(48), 18814–18824.  
1258 <http://doi.org/10.1523/JNEUROSCI.3128-13.2013>

- 1259 Harel, J., Koch, C., & Perona, P. (2007). Graph-Based Visual Saliency. In *Advances in Neural*  
1260 *Information Processing Systems 19* (Vol. 19, pp. 545–552). The MIT Press.  
1261 <https://doi.org/10.7551/mitpress/7503.003.0073>
- 1262 Hayes, T. R., & Henderson, J. M. (2020). Center bias outperforms image salience but not  
1263 semantics in accounting for attention during scene viewing. *Attention, Perception, and*  
1264 *Psychophysics*, 82(3), 985–994. <https://doi.org/10.3758/s13414-019-01849-7>
- 1265 Hayes, T. R., & Henderson, J. M. (2021). Looking for Semantic Similarity: What a Vector-Space  
1266 Model of Semantics Can Tell Us About Attention in Real-World Scenes. *Psychological*  
1267 *Science*, 32(8), 1262–1270. <https://doi.org/10.1177/0956797621994768>
- 1268 Henderson, J. M., & Hayes, T. R. (2017). Meaning-based guidance of attention in scenes as  
1269 revealed by meaning maps. *Nature Human Behaviour*, 1(October).  
1270 <https://doi.org/10.1038/s41562-017-0208-0>
- 1271 Henderson, J. M., Hayes, T. R., Peacock, C. E., & Rehrig, G. (2021). Meaning maps capture the  
1272 density of local semantic features in scenes: A reply to Pedziwiatr, Kümmerer, Wallis,  
1273 Bethge & Teufel (2021). *Cognition*, (January), 104742.  
1274 <https://doi.org/10.1016/j.cognition.2021.104742>
- 1275 Henderson, J. M., Malcolm, G. L., & Schandl, C. (2009). Searching in the dark: Cognitive  
1276 relevance drives attention in real-world scenes. *Psychonomic Bulletin & Review*, 16(5),  
1277 850–856. <https://doi.org/10.3758/PBR.16.5.850>
- 1278 Horga, G., & Abi-Dargham, A. (2019). An integrative framework for perceptual disturbances in  
1279 psychosis. *Nature Reviews Neuroscience*, 20(12), 763–778.  
1280 <https://doi.org/10.1038/s41583-019-0234-1>
- 1281 Hsieh, P.-J. J., Vul, E., & Kanwisher, N. (2010). Recognition alters the spatial pattern of fMRI  
1282 activation in early retinotopic cortex. *Journal of Neurophysiology*, 103(3), 1501–1507.  
1283 <https://doi.org/10.1152/jn.00812.2009>
- 1284 Hwang, A. D., Wang, H.-C., & Pomplun, M. (2011). Semantic guidance of eye movements in



- 1285 real-world scenes. *Vision Research*, 51(10), 1192–1205.  
1286 <https://doi.org/10.1016/j.visres.2011.03.010>
- 1287 Itti, L., & Koch, C. (2000). A saliency-based search mechanism for overt and covert shifts of  
1288 visual attention. *Vision Research*, 40(10–12), 1489–1506. [https://doi.org/10.1016/S0042-](https://doi.org/10.1016/S0042-6989(99)00163-7)  
1289 6989(99)00163-7
- 1290 Itti, L., & Koch, C. (2001). Computational modelling of visual attention. *Nature Reviews*  
1291 *Neuroscience*, 2(3), 194–203. <https://doi.org/10.1038/35058500>
- 1292 Judd, T., Durand, F., Torralba, A., Azam, S., Gilani, S. O., Jeon, M., ... Torralba, A. (2012). A  
1293 Benchmark of Computational Models of Saliency to Predict Human Fixations. In *MIT*  
1294 *Technical Report* <https://doi.org/10.5220/0005678701340142>
- 1295 Kaspar, K., Hloulal, T. M., Kriz, J., Canzler, S., Gameiro, R. R., Krapp, V., & König, P. (2013).  
1296 Emotions' Impact on Viewing Behavior under Natural Conditions. *PLoS ONE*, 8(1).  
1297 <https://doi.org/10.1371/journal.pone.0052737>
- 1298 Kietzmann, T. C., & König, P. (2015). Effects of contextual information and stimulus ambiguity  
1299 on overt visual sampling behavior. *Vision Research*, 110(Part A), 76–86.  
1300 <https://doi.org/10.1016/j.visres.2015.02.023>
- 1301 Kietzmann, Tim C., Geuter, S., & König, P. (2011). Overt Visual Attention as a Causal Factor of  
1302 Perceptual Awareness. *PLoS ONE*, 6(7), e22614.  
1303 <https://doi.org/10.1371/journal.pone.0022614>
- 1304 Kilpeläinen, M., & Georgeson, M. A. (2018). Luminance gradient at object borders  
1305 communicates object location to the human oculomotor system. *Scientific Reports*, 8(1), 1–  
1306 11. <https://doi.org/10.1038/s41598-018-19464-1>
- 1307 Kleiner, M., Brainard, D., Pelli, D., Ingling, A., Murray, R., & Broussard, C. (2007). What's new in  
1308 Psychtoolbox-3? *Perception*, 36(14).
- 1309 Kourtzi, Z., & Connor, C. E. (2011). Neural representations for object perception: Structure,  
1310 category, and adaptive coding. *Annual Review of Neuroscience*, 34, 45–67.

- 1311 <https://doi.org/10.1146/annurev-neuro-060909-153218>
- 1312 Kriegeskorte, N. (2015). Deep Neural Networks: A New Framework for Modeling Biological  
1313 Vision and Brain Information Processing. *Annual Review of Vision Science*, 1(1), 417–446.  
1314 <https://doi.org/10.1146/annurev-vision-082114-035447>
- 1315 Król, M., & Król, M. (2019). The world as we know it and the world as it is: Eye-movement  
1316 patterns reveal decreased use of prior knowledge in individuals with autism. *Autism*  
1317 *Research*, 12(9), 1386–1398. <https://doi.org/10.1002/aur.2133>
- 1318 Kroner, A., Senden, M., Driessens, K., & Goebel, R. (2020). Contextual encoder–decoder  
1319 network for visual saliency prediction. *Neural Networks*, 129, 261–270.  
1320 <https://doi.org/10.1016/j.neunet.2020.05.004>
- 1321 Kümmerer, M., Bylinskii, Z., Judd, T., Borji, A., Itti, L., Durand, F., ... Torralba, A. (2020).  
1322 MIT/Tübingen Saliency Benchmark. Retrieved from <https://saliency.tuebingen.ai/>
- 1323 Kümmerer, M., Wallis, T. S. A., Gatys, L. A., & Bethge, M. (2017). Understanding Low-and  
1324 High-Level Contributions to Fixation Prediction. *Proceedings of the IEEE International*  
1325 *Conference on Computer Vision (ICCV)*, 4799–4808.  
1326 <https://doi.org/10.1109/ICCV.2017.513>
- 1327
- 1328 Wolf, C., & Lappe, M. (2021). Salient objects dominate the central fixation bias when orienting  
1329 toward images. *Journal of Vision*, 21(8), 1–21. <https://doi.org/10.1167/jov.21.8.23>
- 1330 Lengyel, G., Nagy, M., & Fiser, J. (2021). Statistically defined visual chunks engage object-  
1331 based attention. *Nature Communications*, 12(1), 1–12. [https://doi.org/10.1038/s41467-020-](https://doi.org/10.1038/s41467-020-20589-z)  
1332 [20589-z](https://doi.org/10.1038/s41467-020-20589-z)
- 1333 Lengyel, G., Žalalytė, G., Pantelides, A., Ingram, J. N., Fiser, J., Lengyel, M., & Wolpert, D. M.  
1334 (2019). Unimodal statistical learning produces multimodal object-like representations.  
1335 *ELife*, 8, 1–21. <https://doi.org/10.7554/eLife.43942>
- 1336 Liang, H., Gong, X., Chen, M., Yan, Y., Li, W., & Gilbert, C. D. (2017). Interactions between

- 1337 feedback and lateral connections in the primary visual cortex. *Proceedings of the National*  
1338 *Academy of Sciences of the United States of America*, 114(32), 8637–8642.  
1339 <https://doi.org/10.1073/pnas.1706183114>
- 1340 Lyu, M., Choe, K. W., Kardan, O., Kotabe, H. P., Henderson, J. M., & Berman, M. G. (2020).  
1341 Overt attentional correlates of memorability of scene images and their relationships to  
1342 scene semantics. *Journal of Vision*, 20(9), 2. <https://doi.org/10.1167/jov.20.9.2>
- 1343 Malcolm, G. L., Groen, I. I. A., & Baker, C. I. (2016). Making Sense of Real-World Scenes.  
1344 *Trends in Cognitive Sciences*, 20(11), 843–856. <http://doi.org/10.1016/j.tics.2016.09.003>
- 1345 Marr, D., & Nishihara, H. K. (1978). Representation and recognition of the spatial organization of  
1346 three-dimensional shapes. *Proceedings of the Royal Society of London. Series B,*  
1347 *Containing Papers of a Biological Character. Royal Society (Great Britain)*, 200(1140),  
1348 269–294. <https://doi.org/10.1098/rspb.1978.0020>
- 1349 Masciocchi, C. M., Mihalas, S., Parkhurst, D., & Niebur, E. (2009). Everyone knows what is  
1350 interesting: Salient locations which should be fixated. *Journal of Vision*, 9(11).  
1351 <https://doi.org/10.1167/9.11.1>
- 1352 Mooney, C. M. (1957). Age in the development of closure ability in children. *Canadian Journal of*  
1353 *Psychology*, 11(4), 219–226. <https://doi.org/10.1037/h0083717>
- 1354 Morey, R. D. (2008). Confidence Intervals from Normalized Data: A correction to Cousineau  
1355 (2005). *Tutorials in Quantitative Methods for Psychology*, 4(2), 61–64.  
1356 <https://doi.org/10.20982/tqmp.04.2.p061>
- 1357 Morey, R. D., & Rouder, J. N. (2018). *BayesFactor: Computation of Bayes Factors for Common*  
1358 *Designs*. Retrieved from <https://cran.r-project.org/package=BayesFactor>
- 1359 Neri, P. (2014). Semantic control of feature extraction from natural scenes. *Journal of*  
1360 *Neuroscience*, 34(6), 2374–2388. <http://doi.org/10.1523/JNEUROSCI.1755-13.2014>
- 1361 Neri, P. (2017). Object segmentation controls image reconstruction from natural scenes. In  
1362 *PLoS Biology* (Vol. 15). <https://doi.org/10.1371/journal.pbio.1002611>

- 1363 Noton, D., & Stark, L. (1971). Scanpaths in Eye Movements during Pattern Perception. *Science*,  
1364 171(3968), 308–311. <https://doi.org/10.1126/science.171.3968.308>
- 1365 Nuthmann, A., & Henderson, J. M. (2010). Object-based attentional selection in scene viewing.  
1366 *Journal of Vision*, 10(8), 20. <https://doi.org/10.1167/10.8.20>
- 1367 Nuthmann, A., Schütz, I., & Einhäuser, W. (2020). Saliency-based object prioritization during  
1368 active viewing of naturalistic scenes in young and older adults. *Scientific Reports*, 10(1),  
1369 22057. <https://doi.org/10.1038/s41598-020-78203-7>
- 1370 Ongchoco, J. D. K., & Scholl, B. J. (2019). How to Create Objects With Your Mind: From Object-  
1371 Based Attention to Attention-Based Objects. *Psychological Science*, 30(11), 1648–1655.  
1372 <https://doi.org/10.1177/0956797619863072>
- 1373 Pajak, M., & Nuthmann, a. (2013). Object-based saccadic selection during scene perception:  
1374 Evidence from viewing position effects. *Journal of Vision*, 13(2013).  
1375 <https://doi.org/10.1167/13.5.2>
- 1376 Pedziwiatr, M. A., Kümmerer, M., Wallis, T. S. A., Bethge, M., & Teufel, C. (2021a). Meaning  
1377 maps and saliency models based on deep convolutional neural networks are insensitive to  
1378 image meaning when predicting human fixations. *Cognition*, 206(10), 104465.  
1379 <https://doi.org/10.1016/j.cognition.2020.104465>
- 1380 Pedziwiatr, M. A., Kümmerer, M., Wallis, T. S. A., Bethge, M., & Teufel, C. (2021b). There is no  
1381 evidence that meaning maps capture semantic information relevant to gaze guidance:  
1382 Reply to Henderson, Hayes, Peacock, and Rehrig (2021). *Cognition*, (April), 104741.  
1383 <https://doi.org/10.1016/j.cognition.2021.104741>
- 1384 Pilarczyk, J., & Kuniecki, M. J. (2014). Emotional content of an image attracts attention more  
1385 than visually salient features in various signal-to-noise ratio conditions. *Journal of Vision*,  
1386 14(12), 4–4. <https://doi.org/10.1167/14.12.4>
- 1387 Powers, A. R., Mathys, C., & Corlett, P. R. (2017). Pavlovian conditioning–induced  
1388 hallucinations result from overweighting of perceptual priors. *Science*, 357(6351), 596–600.

- 1389 <https://doi.org/10.1126/science.aan3458>
- 1390 R Core Team. (2020). *R: A language and environment for statistical computing*. Vienna: R  
1391 Foundation for Statistical Computing. Retrieved from <https://www.r-project.org/>
- 1392 Rosenholtz, R. (2016). Capabilities and Limitations of Peripheral Vision. *Annual Review of*  
1393 *Vision Science*, 2, 437–457. <https://doi.org/10.1146/annurev-vision-082114-035733>
- 1394 Schütt, H. H., Rothkegel, L. O. M., Trukenbrod, H. A., Engbert, R., & Wichmann, F. A. (2019).  
1395 Disentangling bottom-up versus top-down and low-level versus high-level influences on  
1396 eye movements over time. *Journal of Vision*, 19(3). <https://doi.org/10.1167/19.3.1>
- 1397 Self, M. W., Jeurissen, D., van Ham, A. F., van Vugt, B., Poort, J., & Roelfsema, P. R. (2019).  
1398 The Segmentation of Proto-Objects in the Monkey Primary Visual Cortex. *Current Biology*,  
1399 29(6), 1019-1029.e4. <https://doi.org/10.1016/j.cub.2019.02.016>
- 1400 Self, M. W., van Kerkoerle, T., Supèr, H., & Roelfsema, P. R. (2013). Distinct Roles of the  
1401 Cortical Layers of Area V1 in Figure-Ground Segregation. *Current Biology*, 2121–2129.  
1402 <https://doi.org/10.1016/j.cub.2013.09.013>
- 1403 Stoll, J., Thrun, M., Nuthmann, A., & Einhäuser, W. (2015). Overt attention in natural scenes:  
1404 Objects dominate features. *Vision Research*, 107, 36–48.  
1405 <https://doi.org/10.1016/j.visres.2014.11.006>
- 1406 Tatler, B. W. (2007). The central fixation bias in scene viewing: Selecting an optimal viewing  
1407 position independently of motor biases and image feature distributions. *Journal of Vision*,  
1408 7(14). <https://doi.org/10.1167/7.14.4>
- 1409 Tatler, B. W., & Vincent, B. T. (2009). The prominence of behavioural biases in eye guidance.  
1410 *Visual Cognition*, 17(6–7), 1029–1054. <https://doi.org/10.1080/13506280902764539>
- 1411 Teufel, C., Dakin, S. C., & Fletcher, P. C. (2018). Prior object-knowledge sharpens properties of  
1412 early visual feature-detectors. *Scientific Reports*, 8(1). [https://doi.org/10.1038/s41598-018-](https://doi.org/10.1038/s41598-018-28845-5)  
1413 28845-5
- 1414 Teufel, C., & Fletcher, P. C. (2020). Forms of prediction in the nervous system. *Nature Reviews*

- 1415 *Neuroscience*, 21(4), 231–242. <https://doi.org/10.1038/s41583-020-0275-5>
- 1416 Teufel, C., Subramaniam, N., Dobler, V., Perez, J., Finnemann, J., Mehta, P. R., ... Fletcher, P.  
1417 C. (2015). Shift toward prior knowledge confers a perceptual advantage in early psychosis  
1418 and psychosis-prone healthy individuals. *Proceedings of the National Academy of*  
1419 *Sciences*, 112(43), 13401–13406. <https://doi.org/10.1073/pnas.1503916112>
- 1420 Van der Linden, L., Mathôt, S., & Vitu, F. (2015). The role of object affordances and center of  
1421 gravity in eye movements toward isolated daily-life objects. *Journal of Vision*, 15(5), 1–18.  
1422 <https://doi.org/10.1167/15.5.8>
- 1423 Vincent, B. T., Baddeley, R., Correani, A., Troscianko, T., & Leonards, U. (2009). Do we look at  
1424 lights? Using mixture modelling to distinguish between low- and high-level factors in natural  
1425 image viewing. *Visual Cognition*, 17(6–7), 856–879.  
1426 <https://doi.org/10.1080/13506280902916691>
- 1427 Vö, M. L. H., & Schneider, W. X. (2010). A glimpse is not a glimpse: Differential processing of  
1428 flashed scene previews leads to differential target search benefits. *Visual Cognition*, 18(2),  
1429 171–200. <https://doi.org/10.1080/13506280802547901>
- 1430 Wagemans, J., Elder, J. H., Kubovy, M., Palmer, S. E., Peterson, M. A., Singh, M., & von der  
1431 Heydt, R. (2012). A century of Gestalt psychology in visual perception: I. Perceptual  
1432 grouping and figure-ground organization. *Psychological Bulletin*, 138(6), 1172–1217.  
1433 <https://doi.org/10.1037/a0029333>
- 1434 Wallis, T. S. A., Bethge, M., & Wichmann, F. A. (2016). Testing models of peripheral encoding  
1435 using metamerism in an oddity paradigm. *Journal of Vision*, 16(2).  
1436 <http://doi.org/10.1167/16.2.4>
- 1437 Walsh, K. S., McGovern, D. P., Clark, A., & O'Connell, R. G. (2020). Evaluating the  
1438 neurophysiological evidence for predictive processing as a model of perception. *Annals of*  
1439 *the New York Academy of Sciences*, 1464(1), 242–268. <https://doi.org/10.1111/nyas.14321>
- 1440 Wilming, N., Betz, T., Kietzmann, T. C., & König, P. (2011). Measures and Limits of Models of

- 1441            Fixation Selection. *PLoS ONE*, 6(9). <https://doi.org/10.1371/journal.pone.0024038>
- 1442    Wynn, J. S., Shen, K., & Ryan, J. D. (2019). Eye movements actively reinstate spatiotemporal  
1443            mnemonic content. *Vision*, 3(2). <https://doi.org/10.3390/vision3020021>
- 1444    Zelinsky, G. J., & Bisley, J. W. (2015). The what, where, and why of priority maps and their  
1445            interactions with visual working memory. *Annals of the New York Academy of Sciences*,  
1446            1339(1), 154–164. <https://doi.org/10.1111/nyas.12606>