

# Effortless integration of probabilistic visual input

Andrey Chetverikov<sup>1</sup>, Árni Kristjánsson<sup>2</sup>

1 - Donders Institute for Brain, Cognition and Behavior, Radboud University, The Netherlands,  
[a.chetverikov@donders.ru.nl](mailto:a.chetverikov@donders.ru.nl)

2 - Icelandic Vision Lab, Faculty of Psychology, University of Iceland, Iceland; School of Psychology,  
National Research University Higher School of Economics, Russian Federation, [ak@hi.is](mailto:ak@hi.is)

## Abstract

Prominent theories of perception suggest that the brain builds probabilistic models of the world, assessing the statistics of the visual input to inform this construction. However, the evidence for this idea is often based on simple impoverished stimuli, and the results have often been discarded as an illusion reflecting simple “summary statistics” of visual inputs. Here we show that the visual system represents probabilistic distributions of complex heterogeneous stimuli. Importantly, we show how these statistical representations are integrated with representations of other features and bound to locations, and can therefore serve as building blocks for object and scene processing. We uncover the organization of these representations at different spatial scales by showing how expectations for incoming features are biased by neighboring locations. We also show that there is not only a bias, but also a skew in the representations, arguing against accounts positing that probabilistic representations are discarded in favor of simplified summary statistics (e.g., mean and variance). In sum, our results reveal detailed probabilistic encoding of stimulus distributions, representations that are bound with other features and to particular locations.

## Introduction

How the brain represents the visual world is a long-standing question in cognitive science. One captivating idea is that the brain builds statistical models that describe probability distributions of visual features in the environment<sup>1-7</sup>. By combining information about different features and their locations, the brain can then form representations of objects and scenes. Indeed, the idea that the brain represents feature distributions matches our conscious visual experience well. Most objects, such as the apple in Figure 1A, contain a multitude of feature values that can be quantified as a probability distribution, and we are seemingly aware of these feature constellations. Surprisingly, most studies of probabilistic representations do not test how such constellations are represented, assuming instead that a stimulus is described by a single value, such as the orientation of a Gabor patch in vision studies or the hue of an item in working memory experiments and that the only uncertainty comes from the sensory noise. While this unrealistic assumption was noted<sup>3</sup> early on, it is still prevalent, leaving open the possibility that the results can be explained with alternative models without assuming detailed representations of probability distributions<sup>8-10</sup>.

Here, we aim to close this gap and ask 1) if the visual system is capable of quickly forming precise representations of heterogeneous stimuli, representations that reflect the probability distribution of their features and 2) if such representations can be bound to other features or to spatial locations thereby serving as building blocks for upstream object and scene processing.

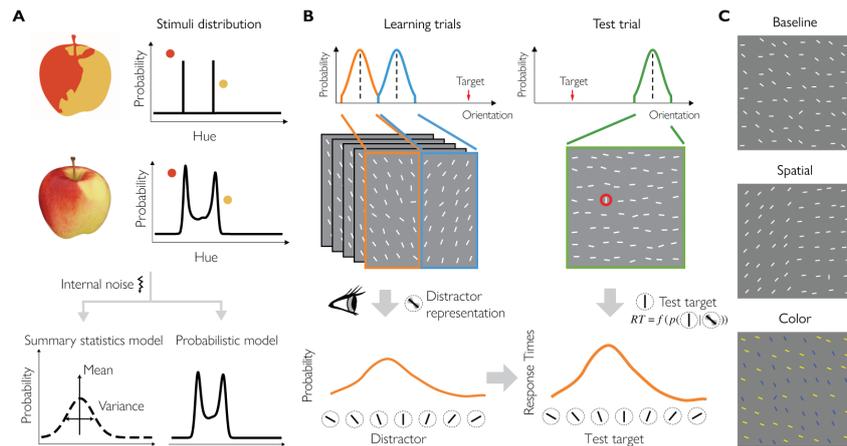


Figure 1. General approach and methods. A: A typical stimulus used to study probabilistic perception involves an impoverished version of the environment, similar to a sketch of an apple (top-left). The hues of this stimulus can be quantified as a discrete probability distribution with only a few probable values (top-right). In contrast, real objects have a multitude of feature values corresponding to a complex-shaped probability distribution (middle). An accurate probabilistic model would maintain the important details of the distribution as much as internal noise permits, while a summary statistics model suggests that probabilities are represented as a combination of simple parameters, such as mean and variance (bottom). B: In our experiments, in each block observers searched for an odd-one-out line among distractors. On learning trials (upper-left), distractors were drawn from two distributions that were either mixed together or separated by location or color with one example of the spatial separation shown here. We assumed that observers would form a distractor representation by learning which distractors are more probable (bottom-left). On test trials (upper-right), we varied the similarity between the target and previously learned distractors. We then measured response times assuming that they should be monotonically related to the probability of a given target being a distractor based on a simplified ideal observer model (bottom-right). C: Example stimuli used on learning trials in Experiment 1.

1 How can the brain represent heterogeneous stimuli, that is, stimuli that have more than one feature  
 2 value? The visual system may track each feature value at each location to form a representation that  
 3 would be identical to the stimulus. However, this would be extremely costly in terms of  
 4 computational resources and unnecessary or even misleading for action because specific feature  
 5 values can vary from one moment to another because of changes in viewpoint, lighting, etc. Another  
 6 possibility is that only a few values, for example, the mean and the variance (“summary statistics”  
 7 <sup>8,10–13</sup>), are represented. But this is also unlikely because multiple stimuli can have the same summary  
 8 statistics while being quite different from each other. More realistically, the brain could follow the  
 9 middle course by approximating feature distributions in the responses of neuronal populations that  
 10 capture the important aspects of stimuli without being too detailed (Figure 1A).

11 Previous studies have indeed shown that the visual system encodes the approximate distribution of  
 12 visual features and uses them in perceptual decision-making <sup>14,15</sup>. However, most of the findings are  
 13 confined to relatively long-term learning of environmental statistics. If feature probability  
 14 distributions are to be useful for everyday visual tasks, such as object recognition or scene  
 15 segmentation, the brain needs to learn feature distributions quickly and effortlessly. Importantly, we  
 16 have previously provided initial evidence that such rapid learning may occur in simple cases by  
 17 studying how human observers learn to ignore distracting stimuli while searching the visual scene <sup>16–</sup>  
 18 <sup>19</sup>. Observers were asked to find an odd-one-out item in a search array where distractor features  
 19 (colors or orientations) were randomly drawn from a given probability distribution for several trials.  
 20 A test trial was then presented with a target of varying similarity to previously learned distractors.  
 21 We found that response times as a function of this similarity parameter followed the shape of the  
 22 previously learned probability distribution, whether it was Gaussian, uniform, skewed, or even  
 23 bimodal. That is, the search was slowed proportionally to how unexpected the target was, based on  
 24 previously learned environmental statistics. This shows that representations of the shape of feature  
 25 probability distributions in the visual input (similar to scene statistics<sup>20,21</sup>) is not limited to long-term  
 26 learning, but can occur rapidly.

1 This previous work was, however, limited to simple scenarios with a single feature distribution  
2 present, while real environments contain multiple objects and scene parts with different features.  
3 Furthermore, knowledge about statistics of a given feature (e.g., orientation) in isolation is not very  
4 useful. Observers need to know *where* in the external world a given feature distribution is and which  
5 other features should be bound with it (related to the “binding” problem<sup>22</sup>) to recognize objects or  
6 segment scenes. Notably, such binding to spatiotopic locations and to other features does not  
7 necessarily require any additional neural machinery, because information about feature  
8 distributions can be readily encoded in neural population responses<sup>2,3,23,24</sup>. Evidence for such  
9 effortless integration of probabilistic visual inputs is, however, still lacking.

10 Ensemble averaging studies testing how observers estimate probabilistic properties of several sets of  
11 stimuli provide some initial support for this hypothesis. It is well known that observers can estimate  
12 the average of a perceptual ensemble, such as the mean orientation of a set of lines<sup>11,25,26</sup>.  
13 Furthermore, they can estimate properties of subsets grouped by location or by other features  
14 although this causes performance detriments<sup>27–32</sup>. However, this approach has only provided  
15 evidence for single-point estimates (the mean) but not for representations of feature probability  
16 distributions. Here, we aim to fill this gap and test how observers encode properties of feature  
17 distributions and associate them with both spatial locations and other features.

## 18 Results

19 In three experiments, observers viewed dressed-down versions of the environment that allowed  
20 precise control of the critical aspects of feature distributions. Observers searched for an unknown  
21 oddball target that differed from other items in orientation and judged whether it was in the upper  
22 or lower half of the stimulus matrix (Figure 1B). Observers did this quickly and accurately despite not  
23 knowing the target or distractor parameters in advance (average response time across experiments  
24 and conditions  $M = 754$  ms,  $SD = 197$ , proportion correct  $M = 0.90$ ,  $SD = 0.04$ ).

25 In all experiments, trials were organized in blocks of intertwined learning and test trials. In each  
26 block, during five to seven learning trials distractor stimuli were drawn randomly from the same  
27 probability distribution. On test trials, we varied the similarity of the current target to non-targets  
28 from preceding trials (Figure 1B). Using this data, we aimed to understand how observers represent  
29 complex heterogeneous stimuli such as visual search distractors.

30 **Bayesian observer model.** How do behavioral responses depend on distractor representations from  
31 previous trials? To answer this question and to reconstruct distractor representations from the  
32 behavioral responses of our observers, we built a Bayesian memory-guided observer model linking  
33 observers’ internal representations of distractors to response times.

34 Our participants located a target among a set of distractors and indicated if it is in the top or the  
35 lower part of the stimuli matrix. On each trial, the experimenter sets the parameters of the target  
36 feature distribution,  $p(s_i | L_T = i)$ , and of the distractor feature distribution,  $p(s_i | L_T \neq i)$ , for each  
37 location  $i = 1 \dots N$  in the stimuli matrix as well as the target location ( $L_T$ ). These parameters are  
38 then used to generate the stimuli at each location ( $s_i$ ). Neither the task parameters nor the stimuli  
39 are known to the observer.

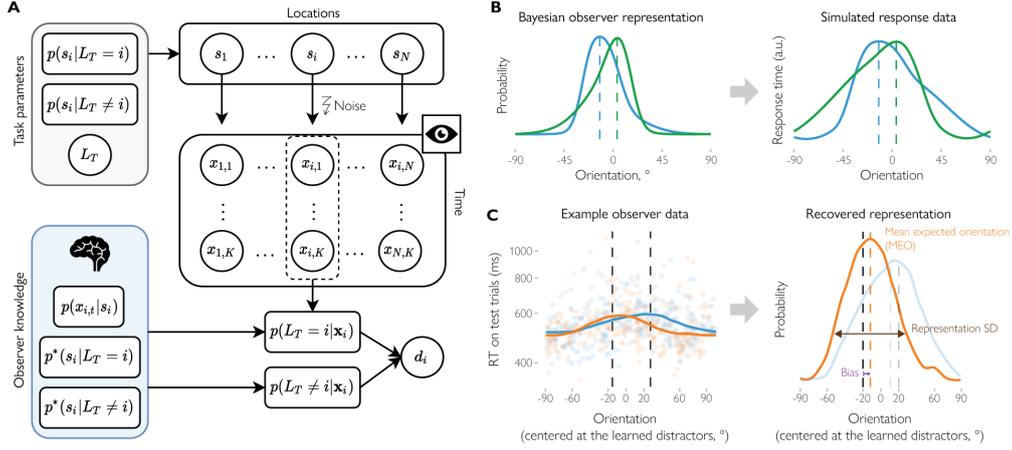


Figure 2. The Bayesian observer model provides a way of reconstructing distractor representations. **A:** The Bayesian observer model. The stimuli  $s_1 \dots s_N$  at different locations are generated on each trial based on task parameters: the target feature distribution  $p(s_i | L_T = i)$ , the distractor feature distribution,  $p(s_i | L_T \neq i)$ , and the target location  $L_T$ . At each moment in time and for each location, observers obtain samples of sensory observations  $x_{i,t}$  corrupted by sensory noise,  $p(x_{i,t} | s_i)$ . Using knowledge about the sensory noise distribution and the approximation of feature distributions for targets and distractors obtained during learning trials,  $p^*(s_i | L_T = i)$  and  $p^*(s_i | L_T \neq i)$ , observers compute probabilities that the sensory observations at a given location correspond to the target,  $p(L_T = i | \mathbf{x}_i)$ , or a distractor,  $p(L_T \neq i | \mathbf{x}_i)$ . These probabilities are combined into a decision variable  $d_i$  used to make a decision or to continue gathering evidence if the currently available observations do not provide enough evidence for the decision (see details in Methods). **B:** The Bayesian observer model enables predictions about response times for a given representation of distractor stimuli (different example distributions are shown in blue and green). Crucially, there is a monotonic relationship between the two, with response times increasing with an increase in distractor probability. **C:** In our analyses, we used the monotonic relationship between probabilistic representations and response times to recover the representation of distractors (right) based on the response times on test trials (left). Here, the data from an example observer in the Spatial condition is split based on whether the target was located in the left (orange) or in the right (blue) hemifield. We then estimated the parameters of the representation, such as the mean expected orientation (dashed orange line), SD and across-distribution bias (the shift in the mean towards the other distribution relative to the true mean, shown by the dashed black line).

- 1 Instead, at each moment in time  $t$ , the observer obtains sensory observations at each location ( $x_{i,t}$ ).
- 2 These observations are not identical to the stimuli because of sensory noise,  $p(x_{i,t} | s_i)$ . In other
- 3 words, a given stimulus might result in different sensory responses, and, conversely, a given sensory
- 4 observation might correspond to different stimuli.
- 5 To find the target, the observer compares for each location the probability that the sensory
- 6 observations are caused by a target present at that location,  $p(L_T = i | \mathbf{x}_i)$  where  $\mathbf{x}_i$  are the samples
- 7 obtained for location  $i$  up until the response time, against the probability that they are caused by a
- 8 distractor,  $p(L_T \neq i | \mathbf{x}_i)$ :

$$9 \quad d_i = \frac{p(L_T = i | \mathbf{x}_i)}{p(L_T \neq i | \mathbf{x}_i)} \quad (1)$$

10 While this decision model is relatively simple, it provides a good intuition for observer behavior in  
 11 the task (a more optimal model is provided in the Supplement but the conclusions do not depend on  
 12 model choice). For this decision rule, the observer representation of distractor features learned from  
 13 previous trials is related to response times:

$$14 \quad RT \approx \frac{C_1}{C_0 - \log p^*(s_i | L_T \neq i, \theta_{prev})} \quad (2)$$

15 where  $C_0$  and  $C_1$  are constants (see details in Methods). In words, there is an inverse relationship  
 16 between response times and the approximate likelihood that a given stimulus is a distractor,  
 17  $p^*(s_i | L_T \neq i, \theta_{prev})$ , with the information obtained from previous trials described by a set of latent  
 18 parameters,  $\theta_{prev}$ . When the probability that a stimulus at a given location (e.g., a test target) is a  
 19 distractor is lower, response times are higher, and vice versa.

1 This model provides an important insight, namely, that observers' representations are monotonically  
2 related to response times (Figure 2B). Hence, the relationship between the distribution parameters  
3 (mean, standard deviation, and skewness) reconstructed from RTs and from the true representation  
4 parameters would hold under any other monotonic transformation (for example, if RTs are log-  
5 transformed and the baseline is subtracted as we do in our analyses; see also Figure S1). In other  
6 words, response times can be used to approximately reconstruct observers' representations of  
7 distractors and estimate their parameters.

8 **Binding orientation probabilities to locations and colors.** Having shown how observer response  
9 times should be related to the distractor representations, we now turn to the empirical data. By  
10 analyzing observers' response times to different test targets, we were able to infer which  
11 orientations were most difficult to find, resulting in the longest response times. Crucially, we were  
12 able to reconstruct observers' representations of the probability distributions that they were  
13 exposed to during learning trials (see Methods).

14 The experiments differed in the structure of the learning trials. There were three conditions in  
15 Experiment 1. The learning trials in the *Spatial* condition were organized so that distractor  
16 distributions in the left and the right hemifield differed to mimic the clustering of similar visual  
17 stimuli in the real world. In the *Color* condition, instead of spatial grouping, different distractor  
18 subsets were grouped by color while individual items were randomly distributed. Finally, in the  
19 *Baseline* condition items from the two distributions had the same color and were randomly  
20 distributed (Figure 1C).

21 Firstly, we report the results on the mean expected orientations (MEO) corresponding to the means  
22 of the recovered representations (Figure 2C). If observers ignore the separation of the two parts of  
23 the distribution, then MEO should match the mean of the overall distribution, but should differ  
24 between the distributions if the representations are bound to locations or colors. For example, if  
25 observers accurately learn the properties of the distributions, the MEO should be at +20° relative to  
26 the overall mean in the Spatial condition when the test line is presented in the hemifield that  
27 previously had distractors with an average relative orientation of +20°.

28 We found that in the Spatial condition, observers' representations in each hemifield followed the  
29 actual physical distractor distribution. The estimated MEO relative to the overall mean was  $M =$   
30  $-14.02^\circ$  ( $SD = 6.02$ ) and  $M = 14.90^\circ$  ( $SD = 5.14$ ) for probes for clockwise (CW) and counterclockwise-  
31 shifted (CCW) distributions, respectively. The difference in MEO between the two distributions was  
32 much larger than zero ( $b = 28.94^\circ$ , 95% HPDI = [25.34, 32.56],  $BF = 6.35 \times 10^{17}$ ) showing that  
33 observers expected different orientations in different hemifields. We then computed the across-  
34 distribution bias by recoding the errors in MEO relative to the true mean for each distribution so that  
35 positive values correspond to shifts towards the other distribution. That is, the bias here represents

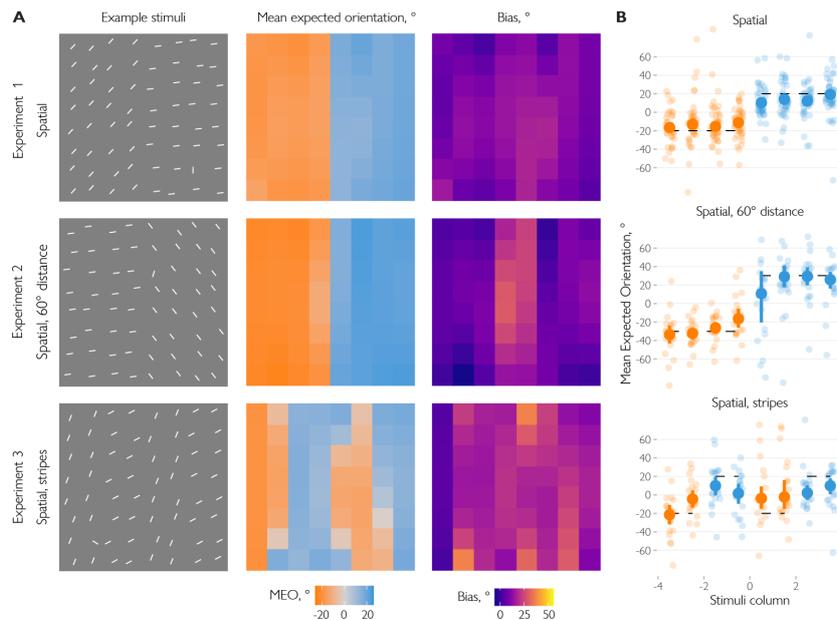


Figure 3. Spatial structure of probabilistic representations. **A:** Example stimuli (left column), recovered mean expected orientations (middle column) and the across-distribution biases in mean expected orientations relative to the true orientations at a given location (right column). The stimuli show a single learning trial from the search task in the corresponding experiment. The mean expected orientation (MEO) was then computed at each location relative to the overall average orientation in the preceding learning block. For presentation purposes, the data were rearranged so that the distribution in the left hemifield (or in the columns 1,2,5,6 in the stripes condition) was oriented clockwise relative to the overall mean. The biases in MEO were computed by subtracting the mean orientation for a given part of the distribution (e.g., at the left/right hemifield in the Spatial condition of Experiment 1) and recoding the resulting errors so that the positive values correspond to a bias towards the other distribution. **B:** Average MEO by column of stimuli matrix in the spatial conditions. Small dots show the data for individual observers, larger dots and bars show means and 95% CI, respectively. Dashed horizontal lines show the true means for a given part of the distribution.

1 by how much observers' expectations deviated from the true mean orientation at a given location  
2 towards the mean orientation at the other location. For both hemifields there was a significant bias  
3 towards the other hemifield ( $M = 5.52^\circ$ , 95% CI = [1.86, 9.14]). This shows that while observers  
4 represent the spatial separation between the two distributions, signals from the other hemifield still  
5 influence their responses.

6 But does spatial separation help observers to track the feature probabilities? In the Baseline  
7 conditions, locations of the CW and CCW distributions were chosen randomly for each learning trial.  
8 We repeated the analysis described above, comparing the response to test targets at the location  
9 that had CW and CCW orientations on immediately preceding trials. We expected to find stronger  
10 across-distribution biases as there was no separation between the distributions across trials.  
11 Importantly, the across-distribution bias was larger in the Baseline (bias  $M = 11.35^\circ$ , 95% CI = [7.71,  
12 15.00]) than the Spatial condition (effect of condition  $M = 5.84$ , 95% CI = [1.10, 10.58],  $BF = 108.24$ ).  
13 In other words, the representations for each distribution were closer to the overall distribution  
14 mean in the Baseline than the Spatial condition. This argues that when the learned distributions are  
15 consistently presented at separate locations, observers can track them better than when they are  
16 mixed.

17 Do observers integrate information about orientation probabilities and color? In the Color condition,  
18 the locations of the test targets were counterbalanced with respect to their colors, so we should  
19 only find differences in MEO if observers formed an association between color and orientation.  
20 Indeed, we found that the MEOs for the two distributions differed ( $b = 7.35$ , 95% HPDI = [1.30,  
21 13.06],  $BF = 148.04$ ) although across-distribution biases were stronger ( $M = 16.30$ , 95% CI = [12.66,  
22 19.86]) than in the Spatial condition ( $M = 10.78$ , 95% CI = [5.99, 15.54],  $BF = 6.56 \times 10^4$ ). This means  
23 that if observers saw yellow lines shifted CW and blue lines shifted CCW relative to the overall

1 distractor mean during learning trials, they learned this association which affected their response  
2 times on subsequent test trials. Importantly, this demonstrates that observers can integrate  
3 information about likely orientations with information about other features (in this case color), even  
4 if there is no spatial information to guide this integration.

5 **Encoding orientation probabilities at different spatial scales.** Having established that observers  
6 associate information about most likely orientations with specific locations or colors, we then asked  
7 if we can uncover the origins of the observed biases by assessing the recovered representations in  
8 the Spatial condition in more detail (for this and later analyses, we increased the sensitivity of our  
9 analyses by combining the data from the Spatial group in Experiment 1 with an additional sample  
10 that performed the same task; see Methods). We computed MEO using the aggregated data from all  
11 participants for each location in the stimuli matrix in this condition. As Figure 3 shows, across-  
12 distribution biases were stronger closer to the boundary between the two hemifields. We then  
13 tested this observation by directly comparing MEOs for test trials with targets presented at the  
14 boundary (two central columns) between the hemifields against other test trials. We found that the  
15 bias was significantly larger at the boundary between the two distributions than in the other  
16 columns ( $M = 4.80^\circ$  ( $SD = 6.99$ ) and  $M = 9.04^\circ$  ( $SD = 11.36$ ),  $b = 4.23$ , 95% HPDI = [0.21, 8.32],  $BF =$   
17  $42.34$ ; Figure 3B). However, the biases were also significantly above zero outside the boundary ( $BF =$   
18  $248$ ). This suggests that the distribution representations are not homogenous and influence each  
19 other strongly when they are close in space, but this mutual influence also extends outside the  
20 immediate neighboring locations (see Discussion).

21 **Bias strength depends on similarity and spatial arrangement.** In two follow-up studies, we further  
22 investigated observers' representations of spatially-grouped heterogeneous stimuli. In Experiment 2,  
23 we tested whether the similarity between the distributions along the tested feature dimension  
24 (orientation) affects the strength of the across-distribution biases. We hypothesized that the bias  
25 should be stronger when the stimuli from the two distributions are more likely to have the same  
26 cause in the external world. For example, the boundary effect in Experiment 1 might occur because  
27 the stimuli close in space are more likely to belong to the same object. By the same reasoning, if the  
28 two distributions are less similar, they are less likely to have the same cause, and the biases should  
29 be weaker.

30 To test this, we used the same spatial arrangement as in the Spatial condition in Experiment 1, but  
31 the distribution means were now  $60^\circ$  away from each other instead of  $40^\circ$  as in Experiment 1 (see  
32 example stimulus in Figure 3A). We found that again, MEO were close to their true values with  $M =$   
33  $26.35^\circ$  ( $SD = 13.43$ ) and  $M = -27.65^\circ$  ( $SD = 10.65$ ) for distributions centred at  $30^\circ$  and  $-30^\circ$  relative to  
34 the overall mean, respectively. Importantly, while there was a strong bias at the boundary between  
35 the distributions,  $M = 19.05^\circ$  ( $SD = 27.27$ ),  $BF = 8.36$ , it was absent at other positions (bias  $M = 0.60^\circ$   
36 ( $SD = 8.65$ ), with  $BF = 4.12$  in favor of no bias). Experiment 2, therefore, shows that reducing the  
37 similarity between the distributions eliminates the biases except for the immediately adjacent  
38 locations.

39 In Experiment 3, we tested whether an even more complex spatial arrangement would allow us to  
40 recover the "map" of observers' expected orientations. To this end, the stimuli were organized in  
41 "stripes" of two matrix columns with two different distributions from Experiment 1 (with means  
42 separated by  $40^\circ$ ) positioned at odd and even stripes (counterbalanced across blocks, Figure 3A). We  
43 found that observers expected clockwise-rotated orientations ( $M = 6.20^\circ$ ,  $SD = 9.91$ ) at locations of

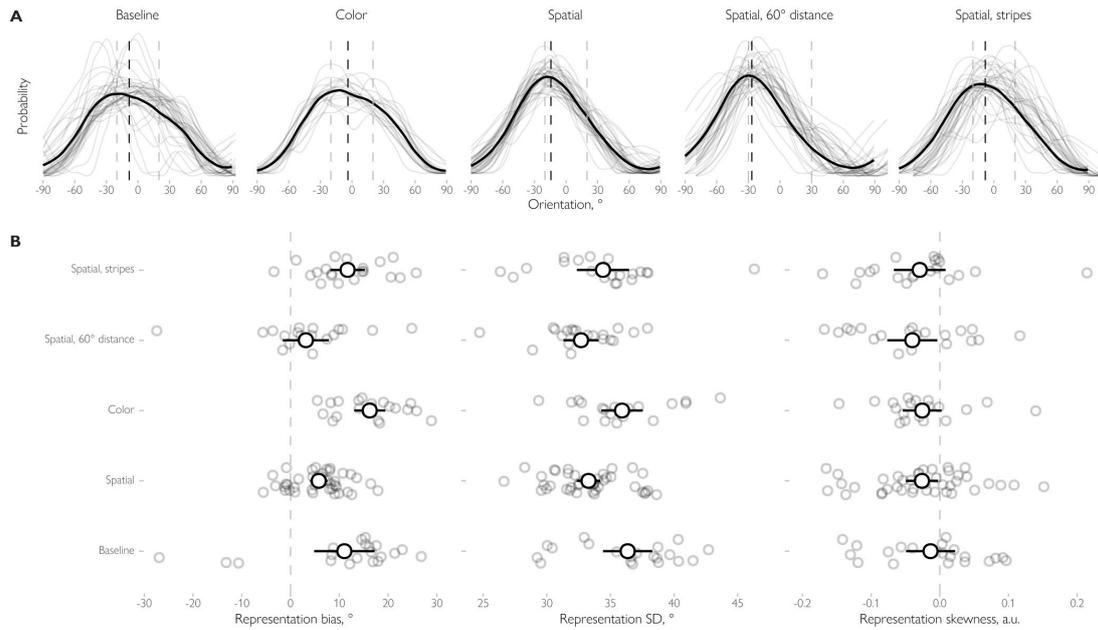


Figure 4. Recovered average representations and their parameters across experiments and conditions. **A:** The black curves show the average representation while representations for individual observers are shown in light gray. Dashed horizontal lines show the mean of the representation (black) and the true mean of the stimulus distributions (light gray). Note that the representations are aligned so that when two distributions are present, the true mean at the tested location is clockwise ( $-20^\circ$  or  $-30^\circ$ ) while the other mean is counterclockwise ( $20^\circ$  or  $30^\circ$  relative to the true mean). **B:** Estimated parameters (bias, SD and skewness). Large dots and errorbars show the mean across observers for a given parameter and the associated 95% confidence intervals. Smaller dots show data for individual subjects.

1 stripes rotated  $20^\circ$  clockwise relative to the overall mean and counterclockwise-rotated orientations  
 2 ( $M = -11.034^\circ$ ,  $SD = 17.11$ ) at other stripe locations. However, the across-distribution bias ( $M =$   
 3  $11.70^\circ$ ,  $SD = 7.52$ ) was stronger than in the Spatial condition in Experiment 1 ( $b = 5.90$ , 95% HPDI =  
 4  $[2.50, 9.33]$ ,  $BF = 4.30$ ). This demonstrates that while separating distributions in space helps  
 5 observers track distributions (as shown in Experiments 1 and 2), the effects of spatial organization  
 6 decrease as the organization becomes more complex.

7 **Higher-order parameters of probabilistic representations.** Next, we asked whether observers'  
 8 representations contain more information about the distributions than just their average? We used  
 9 the reconstructed distractor representations (Figure 4A) to estimate their circular standard deviation  
 10 and circular skewness. The former corresponds to the expected variability among distractors, while  
 11 the latter quantifies their symmetry.

12 First, we hypothesized that if the variability of the distributions is encoded, then the expected  
 13 variability would be higher when distractor distributions are less well separated. Indeed, we found  
 14 that observers' expectations about distractor variability differ between conditions ( $BF = 2.03 \times 10^5$ )  
 15 with lower SD when the distractors were separated by hemifields ( $M = 33.3$ , 95% HPDI =  $[32.2, 34.4]$   
 16 for the Spatial condition with  $40^\circ$  separation and  $M = 32.7$ , 95% HPDI =  $[31.1, 34.2]$  for  $60^\circ$   
 17 separation) compared to other conditions ( $M = 35.9$ , 95% HPDI =  $[34.4, 37.5]$  in the color condition,  
 18  $M = 34.4$ , 95% HPDI =  $[32.9, 35.9]$  for the stripes arrangement condition). When the two  
 19 distributions were less well separated, observers were more uncertain in their estimates, leading to  
 20 distractor representations with higher SD's (Figure 4B).

21 We also expected that the distribution presented at the tested location or in the tested color would  
 22 weigh more highly in the resulting representation, causing an asymmetry. Alternatively, if observers  
 23 only use the mean and variance to encode the distribution (as assumed by "summary statistics"  
 24 accounts), then the represented distribution should be symmetric. We found that observers'

1 representations were asymmetric in all conditions, with a higher probability mass at the side  
2 corresponding to the distribution presented at the tested location or in the tested color,  $M = -0.03$ ,  
3 95% CI = [-0.04, -0.02]. Notably, however, no differences between conditions were found,  $BF =$   
4  $1.99 \times 10^{-6}$ , indicating that symmetry is not affected by the way the distributions are organized in  
5 the display. In sum, observers represent not only the average stimulus values but also their  
6 variability, and the representations are skewed towards distributions presented at other locations or  
7 in different colors.

## 8 Discussion

9 Our main hypothesis was that observers extract information about probabilities of visual features  
10 from heterogeneous stimuli and bind the resulting probabilistic representations with locations on  
11 the one hand and other features on the other. Our results support both these proposals very clearly,  
12 demonstrating how the visual system can build probabilistic representations of the visual world by  
13 extracting information about the features of complex heterogeneous stimuli.

14 A visual search task allowed us to uncover representations of heterogeneous distractors. We  
15 formulated a Bayesian observer model and demonstrated analytically and through simulations that  
16 response times are a monotonic function of observers' expectations about distractor orientations,  
17 supporting earlier empirical findings<sup>16-19</sup>. Using this knowledge, we were able to estimate the  
18 characteristics of observer representations – their means, precision, and skewness – and study how  
19 they vary depending on whether observers can associate them with locations or with other, task-  
20 irrelevant features, such as color.

21 We found that observers encode the feature distributions in scenes containing two different  
22 distributions. The representations generally follow the physical distribution of the stimuli for a given  
23 location or a given color, but importantly, observers are also biased towards the other distribution.  
24 The strength of the bias depends on the degree of separation between the distributions. When the  
25 distributions were separated in space, observers' representations of one distribution were less  
26 influenced by the other distribution, compared to when they were separated by color or were  
27 intermixed (Baseline condition). Furthermore, as we found in Experiment 3, more complex spatial  
28 arrangements ("stripes") increased the biases towards the other distribution. In sum, observers bind  
29 probabilistic representations of visual features to locations and other features, but such binding is  
30 not impenetrable, reminiscent of "illusory conjunctions" of discrete feature values<sup>33</sup>.

31 We were then able to recover the representation of the distribution at different spatial scales. We  
32 found that for spatial separation, the biases are stronger at the boundary between the two  
33 distributions. This is reminiscent of a hierarchical organization of information about feature  
34 probabilities within a scene proposed for perceptual ensembles<sup>11,25</sup>. Such hierarchical ensemble  
35 models suggest that observers represent information about feature probabilities at different levels:  
36 for example, the orientation statistics at a particular location are combined to form a representation  
37 for a group of items, which are, in turn, combined to form an overall ensemble representation. Our  
38 results agree with this idea: the stimuli observers expect at a given location depend not only on what  
39 was previously shown at this location but also on stimuli presented at other locations. Crucially,  
40 biases were also present for the Color condition as well as for the non-boundary locations in the  
41 Spatial condition of Experiment 1. This indicates that the results cannot be explained by purely local  
42 summation of the inputs. It remains to be tested, whether there are actual separable  
43 representations of probability distributions at different levels, or just a unified spatio-featural map  
44 guiding observer responses.

1 We hypothesized that the representations should be more biased by each other when they are more  
2 likely to have the same cause in the external world. This could provide a normative explanation for  
3 the boundary effect: sensory input from adjacent locations is likely to be caused by the same object  
4 and should therefore be integrated while locations far away from each other should be treated  
5 separately. Similarly, for example, in multisensory integration studies, auditory and visual signals are  
6 less likely to be integrated when there is a large discrepancy in their locations<sup>34,35</sup>. However, in  
7 Experiment 1 we found across-distribution biases at locations far from the other distribution. We  
8 reasoned that this is because the stimuli themselves are similar enough to be potentially caused by  
9 the same object, and the inputs are therefore integrated even from non-neighboring locations. In  
10 Experiment 2, we tested this explanation by asking if the similarity between the distributions  
11 themselves in the tested feature domain (orientation) also plays a role. We found that when the  
12 distributions were made more dissimilar, the biases were observed only at the boundary between  
13 the distributions but not at other locations. That is, observers no longer take into account the input  
14 from non-neighboring locations, when stimuli are dissimilar. This supports the proposed normative  
15 explanation and suggests that the principles of information integration for heterogeneous visual  
16 inputs are the same as for other cases, such as multisensory integration or estimation of complex  
17 visual features<sup>35,36</sup>.

18 We then tested if observers represent more than just the mean distractor orientation. We found  
19 that observers represent the distractor variability (i.e., the standard deviation or width of their  
20 representations), which varies in a predictable fashion with the separability between distractor  
21 distributions. When the distractor distributions are poorly separated (e.g., by color only or are  
22 organized in “stripes”), their representations are wider, indicating more uncertainty. Furthermore,  
23 the representations are asymmetric where the tail of the distribution corresponding to the  
24 orientations matching the tested location or color is fatter. That is, observers do not simply  
25 represent the distractors with a (biased) mean and variance, their representations have a complex  
26 shape with more relevant information (e.g., previous orientations at a tested location) weighted  
27 higher and less relevant information (e.g., previous orientations at the other locations) having lower  
28 weight, but still influencing the outcome.

29 These findings indicate that observers represent information about distractor features as a  
30 probability distribution rather than only in terms of the summary statistics, in contrast to popular  
31 ideas of simple “summary statistics”. For example, Treisman<sup>12</sup> argued that statistical processing is a  
32 distinct mode of perceptual and attentional analysis of stimulus sets. She proposed that because of  
33 limited attentional capacity statistical summaries are generated that include the mean, variance, and  
34 perhaps the range. These summaries enable rapid assessment of the general properties and layout  
35 of natural scenes<sup>29,37</sup>. Similarly, Rahnev<sup>10,38</sup> argued that observers represent only a summary  
36 consisting of the most likely stimulus and the associated strength of evidence, and Cohen et al.<sup>8</sup>  
37 used summary statistics to explain the richness of consciousness experience. Our results argue  
38 against such views, since the representations that are bound together are far more detailed than  
39 this implies. That is, the brain might instead approximate the visual input by using a complex set of  
40 parameters to provide accurate descriptions of feature probabilities<sup>39,40</sup>.

41 A recent finding may explain why many previous studies have supported summary statistics  
42 proposals. Hansmann-Roth et al.<sup>41</sup> reasoned that optimal behavior requires the encoding of full  
43 feature distributions, not only summaries, but observers might be unable to explicitly report the full  
44 distribution. This is analogous to how difficult it might be to verbally describe the variety of colors of  
45 an apple without resorting to simplifications (see Figure 1A). Hansmann-Roth et al. tested  
46 observers’ representations both implicitly and explicitly and while explicit judgments were limited to

1 the mean and variance of feature distributions, implicit measures revealed detailed representations  
2 of the same distributions. More information was therefore available to observers than studies of  
3 summary statistics, that have mostly relied on explicit measures, have indicated. Crucially,  
4 Hansmann-Roth et al. were able to uncover why this is: revealing these detailed representations  
5 requires implicit methods, such as we use here.

6 In our experiments, observers learn the distractor feature by combining inputs from heterogeneous  
7 stimuli across several trials in each block, and it can be argued that this is different from perceiving a  
8 single stimulus on a single trial. However, the visual cortex aggregates information at many different  
9 timescales<sup>42</sup>. Even on a single trial, perception unfolds in time and at each moment is dependent on  
10 what has been seen before. And even for a simple stimulus, the visual cortex receives inputs from  
11 many retinal neurons that are affected by processing noise, potentially indistinguishable from the  
12 input from varying features. Indeed, this is why stimulus variability (“external noise”) is often used to  
13 manipulate visual uncertainty<sup>43,44</sup>. We therefore believe that distinguishing “simple” and “complex”  
14 perception is impossible. However, our results clearly show that information about feature  
15 probabilities is available for visually-guided behavior.

16 Taken together, our results show that observers can not only encode probabilities of features from  
17 heterogeneous stimuli in detail but also integrate them with both locations and other features that  
18 have different distributions. These results arguably represent the strongest support yet for the long-  
19 standing idea that the brain builds probabilistic models of the world<sup>1,5–7,24,45,46</sup> and show that  
20 probabilistic representations can serve as building blocks for object and scene processing. Notably,  
21 such representations are not simply limited to summary statistics (e.g., a combination of mean and  
22 variance<sup>8</sup>). Our results also indicate that observers do not represent physical stimuli precisely, but  
23 instead construct an approximation influenced by input from other stimuli. This probabilistic  
24 perspective stands in sharp contrast to views where discrete features of individual stimuli are *either*  
25 bound together to form objects or processed “statistically”<sup>12,40</sup>. Instead, we suggest that the  
26 probabilistic representations are automatically bound to locations and other features since such  
27 binding occurred even though it was not required in the task. Probabilistic representations are  
28 therefore not acquired in isolation but constitute an integral part of perception.

## 29 Methods

30 **Participants.** In total, eighty observers (fifty female, age  $M = 23.10$ ) participated in the experiments.  
31 Twenty observers (ten female, age  $M = 25.45$ ) participated in the first experiment (Baseline, Spatial,  
32 and Color conditions) split across two sessions. Twenty observers (fourteen female, age  $M = 25.00$ )  
33 participated in Experiment 2 (“Spatial, 60° distance”) and another twenty (thirteen female, age  $M =$   
34  $25.45$ ) in Experiment 3 (“Spatial, stripes”). Finally, the data from additional twenty observers  
35 (thirteen female, age  $M = 16.50$ ) were collected for the Spatial condition of Experiment 1 to increase  
36 the sensitivity of the spatial analyses.

37 All were staff or students at the Faculty of Psychology, St. Petersburg State University, Russia, or the  
38 University of Iceland, Iceland. The experiment was approved by local ethics boards and was run in  
39 accordance with the Helsinki declaration. Participants at St. Petersburg State University were  
40 rewarded with 500 rubles (approx. 8 USD) per hour each, participants at the University of Iceland  
41 participated without additional reward. All gave their informed consent before participating. The  
42 participants were naïve to the purposes of the studies. Participants were given ample time for  
43 training until they felt comfortable doing the task (the training time ranged from 5 minutes to one  
44 hour depending on the participant).

1 **Procedure.** In *Experiment 1*, each participant performed a search task in five conditions. In each  
2 condition on each trial, observers were presented with 8×8 matrices of 64 lines (line length: 0.71° of  
3 visual angle; matrix size: 16×16°; uniform noise of ±0.5° was added to each line coordinate). The goal  
4 was to find the odd-one-out line whose orientation differed most from the others. Sessions were  
5 separated into blocks of 5 to 7 learning trials followed by 1 or 2 test trials (the number of trials  
6 chosen randomly for each block; the variation in the number of trials was introduced to decrease the  
7 effect of temporal expectations<sup>47</sup>). During learning trials, the overall mean of distracting items varied  
8 randomly with half of the distractors drawn from one distribution and the other half from another  
9 distribution with the properties of distributions differing between conditions:

10 *Baseline:* two truncated Gaussian distributions with SD = 10° and range of 40°, with means separated  
11 by 40° (±20° relative to the overall mean), all stimuli had the same color (white), positions for each  
12 line within the matrix were chosen randomly.

13 *Spatial:* two distributions (either a truncated Gaussian with SD = 10° and a range of 40° or uniform  
14 with the range of 40°) with means separated by 40° (±20° relative to the overall mean), all stimuli  
15 had the same color (white), one distribution was shown in the left half of the matrix, the other in the  
16 right half.

17 *Color:* the same distributions as in the Spatial condition were used, but lines drawn from one  
18 distribution were blue, while lines from the other distribution were yellow. Positions for each line  
19 within the stimuli matrix were chosen randomly.

20 In all cases, two lines were added to each distractor distribution with their orientation equal to the  
21 minimal and maximal values from that distribution range. As a result, Gaussian and uniform  
22 distributions always had the same range. Target orientation on each trial was drawn randomly from  
23 a uniform distribution ranging between 60° and 120° relative to the mean distractor orientation.

24 On test trials, distractors came from a single Gaussian distribution with SD = 10° (range-restricted in  
25 the same way as described above), while target orientation was determined in the same way as on  
26 the prime trials. In the color condition, half of the lines from that distribution were blue, half were  
27 yellow.

28 The Baseline condition had 2304 trials, while the Spatial and Color conditions had 5376 trials each  
29 with the higher number of trials used in the latter case to counterbalance additional factors  
30 (distribution type combinations).

31 *Experiments 2 and 3* generally followed the same procedure as the Spatial condition of Experiment  
32 1. In Experiment 2 the means of the distributions were separated by 60° (±30° relative to the overall  
33 mean) instead of 40° in Experiment 1. In Experiment 3, the two distributions were separated by 40°,  
34 as in Experiment 1, but arranged in “stripes” so that the lines drawn from the first distribution were  
35 positioned in the 1<sup>st</sup>, 2<sup>nd</sup>, 5<sup>th</sup>, and 6<sup>th</sup> columns of the stimuli matrix while the other columns were  
36 populated with lines from the second distribution.

37 **Data processing.** For our main analyses of interest, incorrect responses were excluded and response  
38 times were log-transformed and centered by subtracting the mean for each participant. Then, to  
39 reduce the noise in RT measurements, spatial and featural confounders were removed. First, the  
40 effect of the distance between target locations on consecutive trials and the effect of the target  
41 location were removed by regressing out the fifth-degree polynomials of the absolute distance (in  
42 degrees of visual angle) between the target locations on the current and the previous trials and the  
43 current targets horizontal and vertical coordinates. Then, we also removed potential influences from

1 the well-known oblique effect (the search speed differs between oblique and cardinal stimuli<sup>45,48</sup> by  
2 regressing out the fifth-degree polynomials of target and distractor obliqueness computed as an  
3 absolute distance in degrees to the nearest cardinal orientation. The regression was run separately  
4 for each experiment and condition.

5 To reconstruct observers' distractor representations, we used the response times on the first test  
6 trial in each block. We then converted response times as a function of the similarity between the  
7 test target and previous distractor mean to a probabilistic representation and estimated its  
8 parameters.

9 To convert the noisy response times into probabilities, we first smoothed RT as a function of the test  
10 target and previous distractor mean using the local regression approach (a generalization of the  
11 moving average) for each observer in each condition. To account for circularity, we appended 1/6 of  
12 the data from each end of the orientation space to the opposite end before smoothing. In analyses  
13 applied to each stimulus location, we further assumed that RTs are a smooth function of the stimuli  
14 matrix row within the local regression while columns of the stimuli matrix were treated  
15 independently. We then transformed a smoothed RT function into a probability mass function by  
16 subtracting the baseline and normalizing to one. Finally, we computed the parameters of the  
17 recovered probabilistic representation: the mean expected orientation (circular mean), circular  
18 standard deviation, and circular skewness as defined by Pewsey<sup>49</sup>. Note that under the hypothesized  
19 Bayesian observer model, the estimated standard deviation and skewness are monotonically related  
20 to the true parameters of the distractor representation but are not identical to it (additionally  
21 confirmed in simulations, Figure S1).

22 **Data analysis.** Unless stated otherwise, we used Bayesian hierarchical regression with *brms*<sup>50</sup>  
23 package in R. Note that while we include Bayes factor values in the description of the results, we  
24 were mostly interested in measuring the effects of the variables of interest in our models, hence the  
25 models included the default flat (uniform) priors for regression coefficients. Given that Bayes factors  
26 are heavily prior-dependent, we believe that the information provided by the 95% highest-density  
27 posterior intervals (HDPI) is more useful for judging the results than the Bayes factors. To make sure  
28 that the conclusions are not dependent on the particular analytic approach, we repeated the  
29 analyses using the conventional frequentist statistical test with the same results (the report using  
30 this approach is provided alongside the data in an online repository, see *Data availability*  
31 statement).

32 **Bayesian observer model.** In our experiments, participants located a target among a set of  
33 distractors and indicated if it is in the top or the lower part of the stimuli matrix. On each trial, the  
34 experimenter sets the task parameters, namely, parameters of the target distribution,  $p(s_i|L_T = i)$ ,  
35 and parameters of the distractor distribution,  $p(s_i|L_T \neq i)$ , for each location  $i = 1 \dots N$  in the stimuli  
36 matrix as well as the target location,  $L_T$ . These parameters were then used to generate the stimuli at  
37 each location,  $s_i$ .

38 Neither the task parameters nor the stimuli are known to the Bayesian observer. Instead, at each  
39 moment in time  $t$ , the observer obtains sensory observations at each location,  $x_{i,t}$ . These  
40 observations are not identical to the stimuli because of the presence of sensory noise,  $p(x_{i,t}|s_i)$ .  
41 That is, a given stimulus might result in different sensory responses, and, conversely, a given sensory  
42 observation might correspond to different stimuli. We assume that the observations are distributed  
43 independently at each location and at each moment in time.

44 To make an optimal decision in a particular task, the observer needs to know the relationship  
45 between the sensory observations and the task-relevant quantities. For the visual search task used

1 in our study, we assumed that observers compare for each location the probability that the sensory  
 2 observations are caused by a target present at that location,  $p(L_T = i | \mathbf{x}_i)$  where  $\mathbf{x}_i =$   
 3  $\{x_{i,1}, x_{i,2}, \dots, x_{i,t=K}\}$  are the samples obtained for location  $i$  up until the time  $K$ , against the  
 4 probability that they are caused by a distractor,  $p(L_T \neq i | \mathbf{x}_i)$ :

$$5 \quad d_i = \frac{p(L_T = i | \mathbf{x}_i)}{p(L_T \neq i | \mathbf{x}_i)} \quad (3)$$

6 The observer then decides that a given item is a target as soon as the decision variable at a given  
 7 location reaches a certain threshold  $B$ . Although this decision rule is not fully optimal, because the  
 8 observer makes a decision for each item individually, it greatly reduces the task complexity, and we  
 9 believe that it allows for a more realistic model (the simulations based on a more complex but more  
 10 optimal model are described in the supplement and lead to identical conclusions).

11 The observer can compute the probability of hypotheses  $L_T = i$  and  $L_T \neq i$  given the sensory data  
 12 using the Bayes rule:

$$13 \quad p(L_T = i | \mathbf{x}_i) = \frac{p(\mathbf{x}_i | L_T = i)p(L_T = i)}{p(\mathbf{x}_i)} \quad (4)$$

14 In words, the probability of a hypothesis that a target is at the given location,  $L_T = i$ , for a set of  
 15 sensory observations  $\mathbf{x}_i$  is equal to the likelihood of the data given this hypothesis multiplied by a  
 16 prior probability for this hypothesis  $p(L_T = i)$  and divided by the probability of the observations  
 17  $p(\mathbf{x}_i)$ .

18 Assuming that the prior  $p(L_T = i) = \frac{1}{N} = 1 - p(L_T \neq i)$  is the same for all locations, the decision  
 19 variable can then be rewritten in log-space as the difference in the log-likelihoods in favor of the two  
 20 hypotheses:

$$21 \quad \log d_i = \sum_{t=1}^K \log(p(x_{i,t} | L_T = i)) - \sum_{t=1}^K \log(p(x_{i,t} | L_T \neq i)) + \log\left(\frac{1}{N-1}\right) \quad (5)$$

22 What are the probabilities of sensory observations under each hypothesis,  $p(x_{i,t} | L_T = i)$  and  
 23  $p(x_{i,t} | L_T \neq i)$ ? To compute them, the observer needs to take into account how the stimuli are  
 24 distributed under each hypothesis and how the sensory noise is distributed for each stimulus. We  
 25 assume that the sensory noise distribution is known for the observer through long-time exposure to  
 26 the visual environment (that is, the observer knows  $p(x_{i,t} | s_i)$ ).

27 However, to determine how probable it is that sensory observations correspond to the search  
 28 target, the observer must also know what defines targets and distractors. The experimenter knows  
 29 that only certain orientations describe a target, but the observer is not omniscient and does not  
 30 know the true distributions of target and distractor stimuli, approximating them instead as  
 31  $p^*(s_i | L_T = i)$  and  $p^*(s_i | L_T \neq i)$ . Then the probability of sensory observations under each  
 32 hypothesis can be computed as:

$$33 \quad p(x_{i,t} | L_T \neq i) = \int p(x_{i,t} | s_i) p^*(s_i | L_T \neq i) ds_i \quad (6)$$

34 The probability distributions  $p^*(s_i | L_T = i)$  and  $p^*(s_i | L_T \neq i)$  correspond to the observer's  
 35 approximate representation of target and distractor distributions. Notably, each of them can be  
 36 further separated into the representation based on the previous trials and the one based on the  
 37 current trial:

$$p^*(s_i|L_T \neq i) \equiv p(s_i|L_T \neq i, \boldsymbol{\theta}) = p(s_i|L_T \neq i, \boldsymbol{\theta}_{prev})p(s_i|L_T \neq i, \boldsymbol{\theta}_{curr}) \quad (7)$$

with  $\boldsymbol{\theta} = \{\boldsymbol{\theta}_{prev}, \boldsymbol{\theta}_{curr}\}$  corresponding to the independent latent variables describing the parameters of the previous and the current trial by the observer (similar equations related to targets are omitted for brevity). In our experiments, by design, the parameters of the current trial are controlled with respect to the current stimuli (i.e., the distractors on the current test trial are drawn from a distribution with a mean from 60° to 120° off the current test target). Hence, only  $p(s_i|L_T \neq i, \boldsymbol{\theta}_{prev})$  matters for relative changes in response times.

In our analyses, we wanted to reconstruct the representation of distractor stimuli using the response times for different test targets. Because the decision time is proportionate to the number of samples when the sampling frequency is constant, we aimed to relate the number of samples  $K$  to an observer's approximate representation of distractors based on the previous trials  $p(s_i|L_T \neq i, \boldsymbol{\theta}_{prev})$ .

Assuming that the sensory observations are obtained with high frequency, we can approximate the total evidence in favor of a given hypothesis:

$$\sum_{t=1}^K \log(p(x_{i,t}|L_T \neq i)) \approx K \left( E \left[ \log(p(x_{i,t}|L_T \neq i)) \right] \right) \quad (8)$$

We expect sensory noise to be low compared to the noise in the target and distractor representations. Then, the following approximation is valid:

$$E \left[ \log(p(x_{i,t}|L_T \neq i)) \right] \propto \log(p^*(s_i|L_T \neq i)) + C \quad (9)$$

where  $C$  is a constant. Similar derivations can be used for the total evidence for the alternative hypothesis  $p(x_{i,t}|L_T = i)$ .

Then, given that a decision is made when  $\log d_i = \log B$ :

$$K = \frac{\log B - \log\left(\frac{1}{N-1}\right)}{E \left[ \log(p(x_{i,t}|L_T = i)) \right] - E \left[ \log(p(x_{i,t}|L_T \neq i)) \right]} \quad (10)$$

Given that the target and distractor parameters are independently manipulated in the experiment,  $E \left[ \log(p(x_{i,t}|L_T = i)) \right]$  can be treated as a constant. Similarly,  $p(s_i|L_T = i, \boldsymbol{\theta}_{curr})$  would be constant as discussed above. Given that  $RT \propto K$ , we can then approximate is as follows:

$$RT \approx \frac{C_1}{C_0 - \log p^*(s_i|L_T \neq i, \boldsymbol{\theta}_{prev})} \quad (11)$$

and

$$\log p^*(s_i|L_T \neq i, \boldsymbol{\theta}_{prev}) = C_0 - C_1 \frac{1}{RT} \quad (12)$$

where  $C_0$  and  $C_1$  are constants. In words, there is an inverse linear relationship between the likelihood that a given stimulus is a distractor (in log-space) and the response times. When this likelihood increases, response times decrease.

This model provides an important insight, namely, that observers' representations are monotonically related to response times. Hence, even though  $C_0$  and  $C_1$  are unknown, the relationship between

1 the moments (mean, standard deviation, and skewness) of observers' representations reconstructed  
2 from RT and the true representations would hold under any other monotonic transformation (for  
3 example, RTs are log-transformed and the baseline RTs are subtracted as we do in our analyses).

4

5

## References

1. Rao, R. P., Olshausen, B. A. & Lewicki, M. S. *Probabilistic models of the brain: Perception and neural function*. (MIT Press, 2002).
2. Pouget, A., Dayan, P. & Zemel, R. S. Information processing with population codes. *Nat. Rev. Neurosci.* **1**, 125–32 (2000).
3. Zemel, R. S., Dayan, P. & Pouget, A. Probabilistic Interpretation of Population Codes. *Neural Comput.* **10**, 403–430 (1998).
4. Lange, R. D., Shivkumar, S., Chattoraj, A. & Haefner, R. M. Bayesian Encoding and Decoding as Distinct Perspectives on Neural Coding. *bioRxiv* 1–16 (2020).
5. Fiser, J., Berkes, P., Orbán, G. & Lengyel, M. Statistically optimal perception and learning: from behavior to neural representations. *Trends Cogn. Sci.* 119–130 (2010) doi:10.1016/j.tics.2010.01.003.
6. Knill, D. C. & Pouget, A. The Bayesian brain: The role of uncertainty in neural coding and computation. *Trends Neurosci.* **27**, 712–719 (2004).
7. Tanrikulu, Ö. D., Chetverikov, A., Hansmann-Roth, S. & Kristjánsson, Á. What kind of empirical evidence is needed for probabilistic mental representations? An example from visual perception. *Cognition* **217**, 104903 (2021).
8. Cohen, M. A., Dennett, D. C. & Kanwisher, N. What is the Bandwidth of Perceptual Experience? *Trends Cogn. Sci.* **20**, 324–335 (2016).
9. Block, N. If perception is probabilistic, why does it not seem probabilistic? *Philos. Trans. R. Soc. B Biol. Sci.* **373**, (2018).
10. Rahnev, D. The case against full probability distributions in perceptual decision making. *bioRxiv* (2017) doi:10.1101/108944.
11. Haberman, J. & Whitney, D. Ensemble perception: Summarizing the scene and broadening the limits of visual processing. in *From perception to consciousness: Searching with Anne Treisman* (eds. Wolfe, J. M. & Robertson, L.) 339–349 (Oxford University Press, 2012).
12. Treisman, A. How the deployment of attention determines what we see. *Vis. Cogn.* **14**, 411–443 (2006).
13. Ariely, D. Seeing sets: Representation by statistical properties. *Psychol. Sci.* **12**, 157–162 (2001).
14. Girshick, A. R., Landy, M. S. & Simoncelli, E. P. Cardinal rules: visual orientation perception reflects knowledge of environmental statistics. *Nat. Neurosci.* **14**, 926–932 (2011).
15. Seriès, P. & Seitz, A. R. Learning what to expect (in visual perception). *Front. Hum. Neurosci.* **7**, 668 (2013).
16. Chetverikov, A., Campana, G. & Kristjánsson, Á. Building ensemble representations: How the shape of preceding distractor distributions affects visual search. *Cognition* **153**, 196–210 (2016).
17. Chetverikov, A., Campana, G. & Kristjánsson, Á. Probabilistic rejection templates in visual working memory. *Cognition* **196**, 104075 (2020).
18. Chetverikov, A., Campana, G. & Kristjánsson, Á. Representing Color Ensembles. *Psychol. Sci.* **28**, 1–8 (2017).
19. Chetverikov, A., Hansmann-Roth, S., Tanrikulu, Ö. D. & Kristjánsson, Á. Feature Distribution Learning (FDL): A New Method for Studying Visual Ensembles Perception with Priming of Attention Shifts. in *Neuromethods* 1–21 (Springer, 2019). doi:10.1007/7657\_2019\_20.
20. Oliva, A. & Torralba, A. Modeling the shape of the scene: A holistic representation of the spatial envelope. *Int. J. Comput. Vis.* **42**, 145–175 (2001).
21. Rosenholtz, R. Capabilities and Limitations of Peripheral Vision. *Annu. Rev. Vis. Sci.* **2**, 437–457 (2016).
22. Treisman, A. The binding problem. *Curr. Opin. Neurobiol.* **6**, 171–178 (1996).
23. Vértés, E. & Sahani, M. Flexible and accurate inference and learning for deep generative models. *Adv. Neural Inf. Process. Syst.* **2018-Decem**, 4166–4175 (2018).
24. Sahani, M. & Dayan, P. Doubly Distributional Population Codes: Simultaneous Representation of Uncertainty and Multiplicity. *Neural Comput.* **15**, 2255–2279 (2003).
25. Alvarez, G. A. Representing multiple objects as an ensemble enhances visual cognition. *Trends Cogn. Sci.* **15**, 122–31 (2011).
26. Whitney, D. & Yamanashi Leib, A. Ensemble Perception. *Annu. Rev. Psychol.* **69**, 105–129 (2018).
27. Attarha, M. & Moore, C. M. The capacity limitations of orientation summary statistics. *Atten. Percept. Psychophys.* **77**, 1116–1131 (2015).
28. Oriet, C. & Brand, J. Size averaging of irrelevant stimuli cannot be prevented. *Vision Res.* **79**, 8–16 (2013).
29. Chong, S. C. & Treisman, A. Statistical processing: computing the average size in perceptual groups. *Vision Res.* **45**, 891–900 (2005).

- 1 30. Attarha, M., Moore, C. M. & Vecera, S. P. Summary statistics of size: Fixed processing capacity for multiple  
2 ensembles but unlimited processing capacity for single ensembles. *J. Exp. Psychol. Hum. Percept. Perform.*  
3 **40**, 1440–9 (2014).
- 4 31. Attarha, M. & Moore, C. M. The perceptual processing capacity of summary statistics between and within  
5 feature dimensions. *J. Vis.* **15**, 9 (2015).
- 6 32. Utochkin, I. S. & Vostrikov, K. O. The numerosity and mean size of multiple objects are perceived  
7 independently and in parallel. *PLoS ONE* **12**, 1–20 (2017).
- 8 33. Treisman, A. & Schmidt, H. Illusory conjunctions in the perception of objects. *Cognit. Psychol.* **14**, 107–141  
9 (1982).
- 10 34. Körding, K. P. *et al.* Causal inference in multisensory perception. *PLoS ONE* **2**, (2007).
- 11 35. Shams, L. & Beierholm, U. R. Causal inference in perception. *Trends Cogn. Sci.* **14**, 425–432 (2010).
- 12 36. Landy, M. S., Banks, M. S. & Knill, D. C. Ideal-Observer Models of Cue Integration. in *Sensory Cue*  
13 *Integration* (eds. Trommershäuser, J., Körding, K. & Landy, M. S.) 5–29 (Oxford University Press, 2011).  
14 doi:10.1093/acprof:oso/9780195387247.003.0001.
- 15 37. Emmanouil, T. A. & Treisman, A. Dividing attention across feature dimensions in statistical processing of  
16 perceptual groups. *Percept. Psychophys.* **70**, 946–954 (2008).
- 17 38. Yeon, J. & Rahnev, D. The suboptimality of perceptual decision making with multiple alternatives. *Nat.*  
18 *Commun.* **11**, 1–12 (2020).
- 19 39. Freeman, J. & Simoncelli, E. P. Metamers of the ventral stream. *Nat. Neurosci.* **14**, 1195–1201 (2011).
- 20 40. Rosenholtz, R. Demystifying visual awareness: Peripheral encoding plus limited decision complexity  
21 resolve the paradox of rich visual experience and curious perceptual failures. *Atten. Percept. Psychophys.*  
22 **82**, 901–925 (2020).
- 23 41. Hansmann-Roth, S., Kristjánsson, Á., Whitney, D. & Chetverikov, A. Dissociating implicit and explicit  
24 ensemble representations reveals the limits of visual perception and the richness of behavior. *Sci. Rep.* 1–  
25 12 (2021) doi:10.1038/s41598-021-83358-y.
- 26 42. de Lange, F. P., Heilbron, M. & Kok, P. How Do Expectations Shape Perception? *Trends Cogn. Sci.* **22**, 764–  
27 779 (2018).
- 28 43. Hénaff, O. J., Boundy-Singer, Z. M., Meding, K., Ziemba, C. M. & Goris, R. L. T. Representation of visual  
29 uncertainty through neural gain variability. *Nat. Commun.* **11**, 1–12 (2020).
- 30 44. Barthelmé, S. & Mamassian, P. Evaluation of Objective Uncertainty in the Visual System. *PLoS Comput.*  
31 *Biol.* **5**, e1000504 (2009).
- 32 45. Chetverikov, A., Campana, G. & Kristjánsson, Á. Learning features in a complex and changing environment:  
33 A distribution-based framework for visual attention and vision in general. in *Progress in Brain Research*  
34 vol. 236 97–120 (Elsevier, 2017).
- 35 46. Orhan, A. E. & Ma, W. J. Neural Population Coding of Multiple Stimuli. *J. Neurosci.* **35**, 3825–3841 (2015).
- 36 47. Shurygina, O., Kristjánsson, Á., Tudge, L. & Chetverikov, A. Expectations and perceptual priming in a visual  
37 search task: Evidence from eye movements and behavior. *J. Exp. Psychol. Hum. Percept. Perform.* **45**, 489–  
38 499 (2019).
- 39 48. Wolfe, J. M., Klempen, N. L. & Shulman, E. P. Which end is up? Two representations of orientation in  
40 visual search. *Vision Res.* **39**, 2075–2086 (1999).
- 41 49. Pewsey, A. The large-sample joint distribution of key circular statistics. *Metrika* **60**, (2004).
- 42 50. Bürkner, P.-C. brms : An R Package for Bayesian Multilevel Models Using Stan. *J. Stat. Softw.* **80**, (2017).

43

## 44 Acknowledgments

45 Supported by a grant from the Icelandic Research Fund (IRF #173947-052) and a grant from the  
46 Russian Foundation for Basic Research (RFBR, #15-36-01358). AC was supported by the Radboud  
47 Excellence Initiative. We are grateful to Alena Begler for the help with data collection and to James  
48 Cooke for his invaluable comments on the manuscript.

## 49 Data availability

50 The data and scripts used for the data analysis in this paper are available from <https://osf.io/5pfyn/>.

1 Supplement 1. Bayesian observer model combining information across locations.

2 The model reported in the main text presents a simplified version of the decision-making process  
 3 assuming that stimuli at each location are analyzed separately. We believe that such a model might  
 4 be more realistic as it greatly simplifies the computations that observers have to make. However, for  
 5 the sake of completeness, here we briefly describe a more complex conditionally-optimal memory-  
 6 guided Bayesian observer model. We refer to this model as conditionally optimal for two reasons.  
 7 First, a memory-guided observer is by definition not fully optimal in our task, where the test trial  
 8 parameters are unrelated to the previous learning trials. However, given that the task parameters  
 9 repeat throughout learning trials, using the information from the previous trials might be beneficial  
 10 when the observer does not know that the trial parameters have changed. Secondly, we assume that  
 11 the observer’s learning or memory about the stimuli features might not be ideal, hence they use the  
 12 approximations of feature distributions. We show that under this more complex and more optimal  
 13 model, the predictions with respect to the monotonic relationship between the response times and  
 14 expected distractor probabilities stay the same.

15 **Task structure.** Participants have to locate a target among a set of distractors and indicate if it is in  
 16 the top or in the lower part of the stimuli matrix. The experimenter sets the task parameters for  
 17 each trial, namely, the target distribution,  $p(s_i|L_T = i)$ , and the distractor distribution,  $p(s_i|L_T \neq i)$ ,  
 18 for each location  $i = 1 \dots N$  in the stimuli matrix (with top half having indices from 1 to  $N/2$  and the  
 19 bottom half from  $\frac{N}{2} + 1$  to  $N$ ) as well as the target location ( $L_T$ ), to generate the stimuli ( $s_i$ ) at each  
 20 location. Here,  $L_T = i$  and  $L_T \neq i$  indicate that the target is or is not at location  $i$ , or in other words,  
 21 that the target location is or is not  $i$ , respectively.

22 **Ideal observer model.** At each moment in time  $t = 1 \dots K$  (with  $K$  as the decision moment) and at  
 23 each location  $i$ , the observer obtains sensory observations  $x_{i,t}$  corrupted by the presence of sensory  
 24 noise:

$$25 \quad p(x_{i,t}|s_i) = f_{VM}(x_i; s_i, \kappa_s)$$

26 where  $f_{VM}$  is a von Mises distribution density with concentration parameter  $\kappa_s$  quantifying the  
 27 amount of noise. We assume that the observations are distributed independently at each location  
 28 and at each moment in time:

$$29 \quad p(\mathbf{X}|\mathbf{s}) = \prod_{i=1}^N p(\mathbf{x}_i|s_i) = \prod_{i=1}^N \prod_{t=1}^K p(x_{i,t}|s_i) \quad (S13)$$

30 To make an optimal decision in a particular task, the observer needs to compare the probability that  
 31 a target is located in the upper half of the stimuli matrix with a probability that it is located in the  
 32 lower half:

$$33 \quad d = \frac{p(C = 1|\mathbf{X})}{p(C = 2|\mathbf{X})} \quad (S14)$$

34 where  $C = 1$  and  $C = 2$  correspond to the two hypotheses about the target location. After applying  
 35 the log transformation, the decision variable can be expressed as a difference in the amount of  
 36 evidence for the two hypotheses:

$$37 \quad \log d = \log p(C = 1|\mathbf{X}) - \log p(C = 2|\mathbf{X}) \quad (S15)$$

1 The decision time assuming a certain threshold  $B$  can then be found as a time  $K$  when the decision  
 2 variable reaches the threshold. The average decision time can be found by estimating when the  
 3 expectation of  $\log d$  becomes equal to  $\log B$ :

$$4 \quad K = \frac{\log B}{E[\log p(C = 1|\mathbf{X})] - E[\log p(C = 2|\mathbf{X})]} \quad (S16)$$

6 The probabilities for each hypothesis  $C = 1$  and  $C = 2$  can be found using the Bayes rule. For  
 7 example, for  $C = 1$ :

$$8 \quad p(C = 1|\mathbf{X}) = \frac{p(\mathbf{X}|C = 1)p(C = 1)}{p(\mathbf{X})} \quad (S17)$$

9 Because the observer does not know what stimuli are presented and only knows the sensory  
 10 observations, the likelihood  $p(\mathbf{x}|C = 1)$  needs to be computed by averaging (marginalizing) over the  
 11 unknown stimuli values:

$$12 \quad p(\mathbf{X}|C = 1) = \int p(\mathbf{X}|\mathbf{s})p(\mathbf{s}|C = 1)d\mathbf{s} \quad (S18)$$

13 Because the target can be only present at one location, the likelihood  $p(\mathbf{x}|C = 1)$  is computed by  
 14 summing over the possibilities of finding a target at each particular location:

$$15 \quad p(\mathbf{X}|C = 1) = \sum_{i=1}^{\frac{N}{2}} \int p(\mathbf{X}|\mathbf{s})p^*(\mathbf{s}|L_T = i, \boldsymbol{\theta})d\mathbf{s} \quad (S19)$$

16 where similarly to the main text, we use an asterisk to denote probability distributions as  
 17 approximated by the observer through a set of parameters related to previous and current trials  $\boldsymbol{\theta} =$   
 18  $\{\boldsymbol{\theta}_{prev}, \boldsymbol{\theta}_{curr}\}$ . That is, we assume that the observer is unaware of the true distributions  $p(s_i|L_T =$   
 19  $i)$  and  $p(s_i|L_T \neq i)$  and approximates them instead using the information available.

20 If a target is at location  $i$ , it cannot be anywhere else. Hence:

$$21 \quad p^*(\mathbf{s}|L_T = i, \boldsymbol{\theta}) = p^*(s_i|L_T = i, \boldsymbol{\theta}) \prod_{j \neq i}^N p^*(s_j|L_T \neq j, \boldsymbol{\theta}) \quad (S20)$$

22 Using Eq. S20, it can be further shown that:

$$23 \quad \int p(\mathbf{X}|\mathbf{s})p^*(\mathbf{s}|L_T = i, \boldsymbol{\theta})d\mathbf{s} = \left[ \prod_j^N \int p(\mathbf{x}_j|s_j)p^*(s_j|L_T \neq j, \boldsymbol{\theta})ds_j \right] \frac{\int p(\mathbf{x}_i|s_i)p^*(s_i|L_T = i, \boldsymbol{\theta})ds_i}{\int p(\mathbf{x}_i|s_i)p^*(s_i|L_T \neq i, \boldsymbol{\theta})ds_i} \quad (S21)$$

24 Note that the product in the square brackets is the same for all locations, and the remaining part of  
 25 the equation is a ratio of the probability that the measurements at a given location are from the  
 26 target against the probability that they are from the distractor, similarly to the model described in  
 27 the main text.

28 The probability that a given stimulus is a target (or a distractor) depends on both the previous and  
 29 the current trial:

$$30 \quad p^*(s_i|L_T = i, \boldsymbol{\theta}) = p^*(s_i|L_T = i, \boldsymbol{\theta}_{prev})p^*(s_i|L_T = i, \boldsymbol{\theta}_{curr}) \quad (S22)$$

31 For each location and each location-specific hypothesis  $L_T = i$  and  $L_T \neq i$ , the current trial  
 32 parameters need to be computed separately because of the nature of the odd-one-out task. A target  
 33 is defined as the item most different from the distractors. For simplicity, we assumed that observers

1 use the following circular normal approximation for the distractors at the current trial based on the  
2 sensory observations:

$$3 \quad p^*(s_i | L_T \neq i, \boldsymbol{\theta}_{curr}) = f_{VM}(s_i; \hat{\mu}_{j \neq i}, \hat{\kappa}_{j \neq i}) \quad (S23)$$

4 In words, when the observer needs to estimate, how likely it is that the stimulus at location  $i$  is a  
5 distractor, the observer approximates the distribution of stimuli as a von Mises (circular normal)  
6 distribution based on the sensory observations from other locations.

7 The observer might use the knowledge that the target distribution in the task design is on average  
8  $90^\circ$  away from the mean of distractors. We again assume a von Mises approximation:

$$9 \quad p^*(s_i | L_T = i, \boldsymbol{\theta}_{curr}) = f_{VM}(s_i; \hat{\mu}_{j \neq i} + 90^\circ, \kappa_T) \quad (S24)$$

10 where  $\kappa_T$  is the expected precision of the target distribution. In contrast to the distractor  
11 distribution precision that could be guessed based on the samples on the current trial ( $\hat{\kappa}_{j \neq i}$ ), the  
12 target distribution precision cannot be estimated on a single trial (there is only one target stimulus in  
13 a given trial) and has to be based on the other sources of information (e.g., learning throughout the  
14 experiment).

15 Given that the measurement noise is independent across locations, the likelihood of the hypothesis  
16  $C = 1$  can be further expressed as:

$$17 \quad p(\mathbf{X} | C = 1) = \left[ \prod_{j=1}^N \int (\mathbf{x}_j | s_j) p^*(s_j | L_T \neq j, \boldsymbol{\theta}) ds_j \right] \sum_{i=1}^N \frac{\int p(\mathbf{x}_i | s_i) p^*(s_i | L_T = i, \boldsymbol{\theta}) ds_i}{\int p(\mathbf{x}_i | s_i) p^*(s_i | L_T \neq i, \boldsymbol{\theta}) ds_i} \quad (S25)$$

18 Then, assuming that the prior probability of each decision alternative is the same, the decision  
19 variable can be expressed in log-space as:

$$20 \quad \log d = \log \left( \sum_{i=1}^N \frac{\int p(\mathbf{x}_i | s_i) p^*(s_i | L_T = i, \boldsymbol{\theta}) ds_i}{\int p(\mathbf{x}_i | s_i) p^*(s_i | L_T \neq i, \boldsymbol{\theta}) ds_i} \right) - \log \left( \sum_{i=\frac{N}{2}+1}^N \frac{\int p(\mathbf{x}_i | s_i) p^*(s_i | L_T = i, \boldsymbol{\theta}) ds_i}{\int p(\mathbf{x}_i | s_i) p^*(s_i | L_T \neq i, \boldsymbol{\theta}) ds_i} \right) \quad (S26)$$

21 The decision time assuming a certain threshold  $B$  can then be found as a time  $K$  when the decision  
22 variable reaches the threshold.

23 **Simulations.** To estimate the behavior of the observer using this model, we simulated the decision-  
24 making process and estimated the mean response times while varying the properties of the  
25 distractor representation  $p^*(s_i | L_T \neq i, \boldsymbol{\theta}_{prev})$ . The task parameters were based on the actual  
26 experiment design. We used 36 stimuli for each trial with one stimulus being the test target ( $s_{L_T}$ )  
27 and the rest being the distractors. The distractors on each simulated trial were distributed as  
28  $p(s_i | L_T \neq i) = f_{VM}(s_i; \mu_D, \kappa_D)$  where  $\mu_D \sim U(s_{L_T} + 60^\circ; s_{L_T} + 120^\circ)$  (that is, the mean of  
29 distractors is set to  $60^\circ$  to  $120^\circ$  away from the test stimulus) and  $\kappa_D = 8.7$  (approximately equivalent  
30 to the standard deviation of  $10^\circ$  in orientation space). The sensory observations were assumed to be  
31 noisy ( $\kappa_s = 2$ , approximately equivalent to the standard deviation of  $24^\circ$  in orientation space; note  
32 that this is the noise level for samples collected at each moment in time). The observers' target  
33 representation was assumed to be linked with to the distractor representation as

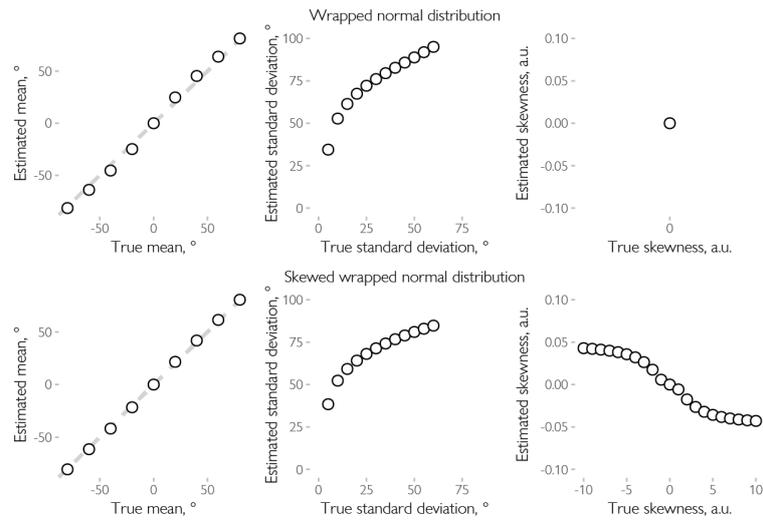
34  $p^*(s_i | L_T = i, \boldsymbol{\theta}_{prev}) = f_{VM}(s_i; \mu_{D,prev}, \kappa_T)$  with  $\kappa_T = 3.35$  (based on a normal approximation to a  
35 uniform target distribution with  $60^\circ$  range used in the experiments). The same  $\kappa_T$  was used for  
36 target-related computations based on the current trial data (Eq. S24). The decision threshold was  
37 set to  $\log B = 4.60$  assuming a 1% probability of error if the observer assumptions are correct. For

1 each test target from  $1^\circ$  to  $180^\circ$  in half-degree steps, we simulated 56 trials for each combination of  
2 distractor representation parameters.

3 We ran simulations for the wrapped skewed normal distribution with the mean varied from  $-60^\circ$  to  
4  $60^\circ$  in  $20^\circ$  steps, while the standard deviation varied from  $20^\circ$  to  $60^\circ$  in  $10^\circ$  steps, and skew varied  
5 from  $-10$  to  $10$  in steps of 2. The results of the simulations (Figure S2) confirmed the findings  
6 obtained with a simplified model: the means are recovered precisely while for standard deviation  
7 and skewness the monotonic relationship holds.

8

1



2

Figure S1. Simulated parameters under the simplified Bayesian observer model. We simulated the response times under the assumptions of the simplified Bayesian observer model described in the main text and applied the same approach as used for the real data to see if the assumed monotonic relationship between the true parameters and the recovered parameters holds. Firstly, we used a simple wrapped normal (top) with means varying from  $-80^\circ$  to  $80^\circ$  in  $20^\circ$  steps and standard deviation from  $5^\circ$  to  $60^\circ$  in  $5^\circ$  steps. For each parameter combination the RT were computed using Eq. 2. We then estimated the parameters of the recovered distribution. As is evident from the plots, the mean estimates were identical to the true mean while the standard deviation was overestimated but the overall monotonic relationship held. The skewness estimate was at zero as expected for the symmetric wrapped normal distribution. Secondly, we simulated the data using the skewed normal distribution (Pewsey, 2008) with means again varying from  $-80^\circ$  to  $80^\circ$  in  $20^\circ$  steps, scale parameter varying from  $5^\circ$  to  $60^\circ$  in  $5^\circ$  steps, and skewness parameter varying from  $-10$  to  $10$  in steps of  $1$ . For the means and standard deviations, the conclusions were the same as for the wrapped normal distribution. Similarly, skewness estimates followed monotonically the changes in the true skewness parameter (note that the sign of the estimated circular skewness is the opposite of the skewness parameter of the skewed wrapped normal distribution because of how it is defined, see Pewsey, 2004). In sum, the mean estimates match the true means, and the standard deviation and skewness estimates monotonically depend on the true standard deviation and the skewness parameters.

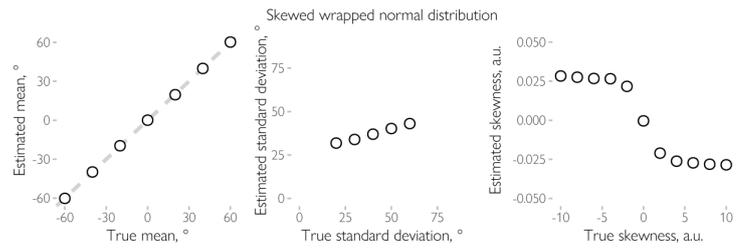


Figure S2. Simulated parameters under the more optimal Bayesian observer model. We simulated the response times under the assumptions of the more complex Bayesian observer model described in the Supplement applied the same approach as used for the real data to see if the assumed monotonic relationship between the true parameters and the recovered parameters holds. The results were similar to the simulations with the simplified model (Figure S1). The mean estimates were identical to the true mean, while for the standard deviation and skewness the monotonic relation holds (note that the sign of the estimated circular skewness is the opposite of the skewness parameter of the skewed wrapped normal distribution because of how it is defined, see Pewsey, 2004). In sum, the mean estimates match the true means, and the standard deviation and skewness estimates monotonically depend on the true standard deviation and the skewness parameters.