NB}
DT}

**ORIGINAL ARTICLE**

# Probabilistic representations as building blocks for higher-level vision

## Andrey Chetverikov[1]    |    Árni Kristjánsson[2]

[1]Donders Institute for Brain, Cognition and Behavior, Radboud University, The Netherlands

[2]Icelandic Vision Lab, Faculty of Psychology, University of Iceland, Iceland

**Correspondence**
Andrey Chetverikov, Donders Institute for Brain, Cognition and Behavior, Radboud University, The Netherlands; Árni Kristjánsson, Faculty of Psychology, University of Iceland, Iceland
Email: a.chetverikov@donders.ru.nl; ak@hi.is

Current theories of perception suggest that the brain represents features of the world as probability distributions, but can such uncertain foundations provide the basis for everyday vision? Perceiving objects and scenes requires knowing not just how features (e.g., colors) are distributed but also where they are and which other features they are combined with. Using a Bayesian computational model, we recovered probabilistic representations used by human observers to search for odd stimuli among distractors. Importantly, we found that the brain integrates information between feature dimensions and spatial locations, leading to more precise representations compared to when information integration is not possible. We also uncovered representational asymmetries and biases, showing their spatial organization and explain how this structure argues against "summary statistics" accounts of visual representations. Our results confirm that probabilistically encoded visual features are bound with other features and to particular locations, providing a powerful demonstration of how probabilistic representations can be a foundation for higher-level vision.

**KEYWORDS**
probabilistic perception, binding problem, ensemble perception, summary statistics, visual search

---

**Abbreviations:** FDL, feature distribution learning; MEO, mean expected orientation.

## 1 | 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 (Fiser et al., 2010; Knill & Pouget, 2004; Lange et al., 2020; Pouget et al., 2000; Rao et al., 2002; Tanrıkulu, Chetverikov, Hansmann-Roth, et al., 2021; Zemel et al., 1998). 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 a while ago (Zemel et al., 1998), it is still prevalent, leaving open the possibility that the results can be explained with alternative models without assuming detailed representations of probability distributions (Block, 2018; Rahnev, 2017).

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.

What precisely do we mean by probabilistic perceptual representations? We assume that the brain operates with probabilistic representations if any of the internal variables used in the perceptual decision-making is represented probabilistically, that is, allowing for uncertainty in their values (similar to, e.g., Koblinger et al., 2021). Often, the concept of probabilistic representations is embedded in the context of Bayesian models. Bayesian perceptual models assume that observers are presented with a stimulus *s* that generates sensory observations *x* with a certain probability *p(x|s)*. The observer knows the parameters of this *generative model*, that is, *p(x|s)*, and can inverse it to compute *p(s|x)*, the probability that a distal stimulus *s* has a certain value given a sensory observation *x*. Importantly, this approach focuses on how a stimulus is inferred from sensory observations. This inferred probability *p(s|x)* is associated with the probabilistic representation. This is intuitively agreeable as long as the focus is on a simple single stimulus but becomes murky when stimuli are heterogenous like the ones shown in Figure 1 (and in reality, homogeneous stimuli do not exist) as it is not clear what an observer infers or should infer in this case.

We approach the problem of probabilistic representations differently, asking instead whether an observer can represent a probability that a stimulus (e.g., an apple or a set of lines in Figure 1) could have a feature with a certain value (e.g., the red color on an apple or a line with a certain orientation). For example, are apples likely to contain red? In Bayesian terms, this means extending the model outlined above to a stimulus-feature-observation hierarchy and asking whether observers can represent probability distributions within this hierarchy, specifically, whether a probability distribution of features given a heterogeneous stimulus is represented within an observers' generative model.

### 1.1 | Probabilistic representations of heterogenous stimuli

How can the brain represent heterogeneous stimuli, that is, stimuli that have more than one feature value? The visual system may track each feature value at each location to form a precise representation isomorphic to the stimulus. However, this would be extremely costly in terms of computational resources and unnecessary or even misleading

for action because specific feature values can vary from one moment to another because of changes in viewpoint, lighting, etc. (Kristjánsson, 2022). Another possibility, explored in the "summary statistics" [1] or "ensemble perception" literature (Ariely, 2001; Cohen et al., 2016; Haberman & Whitney, 2012; Rahnev, 2017; Treisman, 2006) is that only a few values, for example, the mean and the variance are represented. Note that the concept of probabilistic representations is rarely used in this field, because an observer can compute summary statistics without using probability distributions. For example, an average of the features can be computed arithmetically from sensory observations. However, such representations are functionally equivalent to a simplified probability distribution (Figure 1A). But we believe that such simplified representations are also unlikely because multiple stimuli can have the same summary statistics while being quite different from each other. More realistically, the brain could compromise by approximating feature distributions in the responses of neuronal populations that capture important aspects of stimuli without being too detailed (Figure 1A).

Previous studies have indeed shown that the visual system encodes the approximate distribution of visual features and uses them in perceptual decision-making (Girshick et al., 2011; Seriès & Seitz, 2013). However, most of the findings are confined to relatively long-term learning of environmental statistics. If feature probability distributions are to be useful for everyday visual tasks, such as object recognition or scene segmentation, the brain needs to learn feature distributions quickly and effortlessly. Importantly, we have recently provided evidence that such rapid learning may occur in simple cases by studying how human observers learn to ignore distracting stimuli while searching the visual scene (Chetverikov et al., 2016, 2017d, 2020; Chetverikov et al., 2019; Hansmann-Roth et al., 2019; Tanrıkulu, Chetverikov, & Kristjánsson, 2021). The basic idea with this *Feature Distribution Learning* paradigm is to use role-reversal effects upon response times when targets and distractors change their roles between visual search trials, to reveal observers expectations about upcoming search displays (Chetverikov et al., 2019). Priming in visual search is a well-known phenomenon: search is faster when features of targets or distractors repeat from trial-to-trial even when observers do not have to rely on previous trials, that is, when a target is defined as an odd-one-out (Kristjánsson, 2022; Kristjánsson & Campana, 2010; Kristjánsson & Driver, 2008; Lamy et al., 2008; Maljkovic & Nakayama, 1994). And when the targets and distractor switch features ("role reversal"), the search is slower. In our previous experiments, observers were asked to find an odd-one-out item in a search array where, importantly, distractor features (colors or orientations) are randomly drawn from a given probability distribution for several trials rather than having constant features. A test trial (introducing the role reversal) is then presented with a target of varying similarity to previously learned distractors. We found that response times as a function of this similarity parameter followed the shape of the previously learned probability distribution, whether it was Gaussian, uniform, skewed, or even bimodal. That is, the search was slowed proportionally to how unexpected the target was, based on previously learned environmental statistics. This shows that representations of the shape of feature probability distributions in the visual input (similar to scene statistics (Oliva & Torralba, 2001; Rosenholtz, 2016)) are not limited to long-term learning, but can occur rapidly.

This previous work was, however, limited to simple scenarios with a single feature distribution present, while real environments contain multiple objects (that contain multiple features) and scene parts with various different features. Furthermore, knowledge about statistics of a given feature (e.g., orientation) in isolation is not very useful. Observers need to know *where* in the external world a given feature distribution is and which other features should be bound with it (related to the "binding" problem, (Treisman, 1996)) to recognize objects or segment scenes. Notably, such

---

[1] Note that this is different from image-computable summary statistics approaches based on the statistics of the outputs from multi-level image processing filters (Balas et al., 2009; Freeman & Simoncelli, 2011; Portilla & Simoncelli, 2000). While these are related, the statistics in the ensemble perception literature are conceptualized in a more abstract way, more consistent with the type of questions we are interested in here. Yet, even in the image-computable statistics literature it has been demonstrated that images identical in a model statistical space might be still distinguishable by the observers, suggesting that the image-computable summary statistics do not fully match human perception (Wallis et al., 2016)
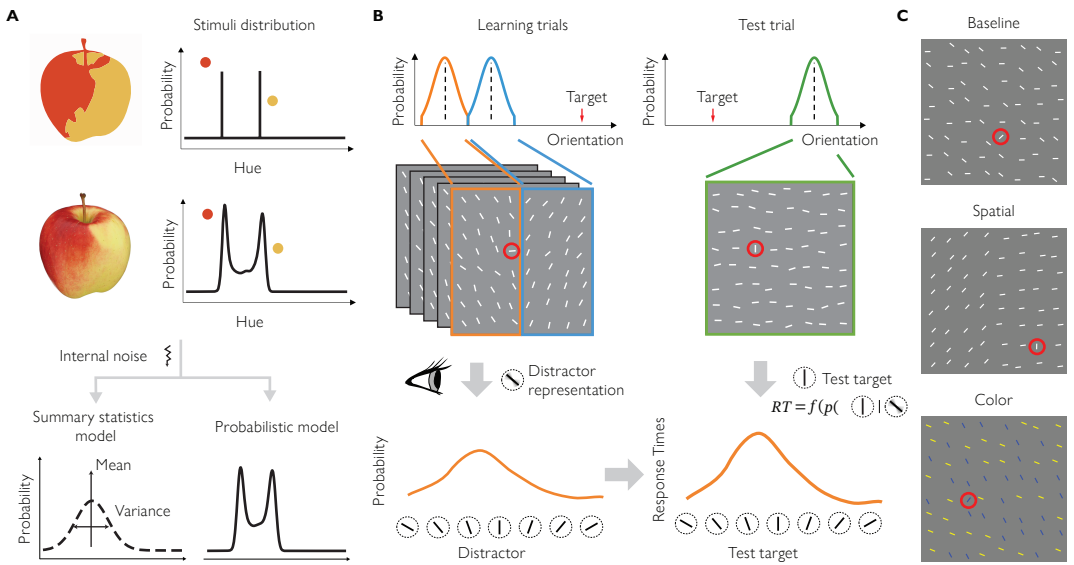
**FIGURE 1** General approach and methods. **A**: A typical stimulus used to study probabilistic perception involves an impoverished version of the environment, similar to a sketch of an apple (top-left). The hues of this stimulus can be quantified as a discrete probability distribution with only a few 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 miniblock (consisting of a sequence of learning trials and 1 or 2 test trials) 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 as shown in previous studies (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.

binding to spatiotopic locations and to other features does not necessarily require any additional neural machinery, because information about feature distributions can be readily encoded in neural population responses (Pouget et al., 2000; Sahani & Dayan, 2003; Vértes & Sahani, 2018; Zemel et al., 1998). Evidence for such effortless integration of probabilistic visual inputs is, however, still lacking.

Ensemble averaging studies testing how observers estimate probabilistic properties of several sets of stimuli provide some initial support for the hypothesis that probabilistic information can be bound to locations or other features. It is well known that observers can estimate the average of a perceptual ensemble, such as the mean orientation of a set of lines (Alvarez, 2011; Haberman & Whitney, 2012; Whitney & Yamanashi Leib, 2018). Notably, they can estimate properties of subsets grouped by location or by other features although this causes performance detriments (Attarha & Moore, 2015a, 2015b; Attarha et al., 2014; Chong & Treisman, 2005; Oriet & Brand, 2013; Utochkin & Vostrikov, 2017). This means that at least a summary representation, functionally equivalent to a simplified probabilistic repre-

sentation based on mean and variance, can be bound to a location. If, for example, a mean can be computed only for a whole set, separate probabilistic representations of different subsets would, in our opinion, be less likely. Yet, this approach has only provided evidence for single-point estimates (e.g., the mean) but no direct evidence for binding of feature probability distributions. Here, we aim to overcome the limitations of previous studies and test how observers encode properties of feature distributions and bind them with both spatial locations and other features.

## 2 | RESULTS

In three experiments, observers viewed dressed-down versions of the environment that allowed precise control of the critical aspects of feature distributions. Observers searched for an unknown oddball target that differed from other items in orientation and judged whether it was in the upper or lower half of the stimulus matrix (Figure 1B). Observers did this quickly and accurately despite not knowing the target or distractor parameters before each block (average response time across experiments and conditions $M$ = 754 ms, $SD$ = 197, proportion correct $M$ = 0.90, $SD$ = 0.04; see Figure S1 for raw RT on test trials by condition).

In all experiments, the trials were organized in miniblocks of intertwined learning and test trials. In each miniblock, during five to seven learning trials, distractor stimuli were randomly drawn from two probability distributions, that were the same within each miniblock but different between miniblocks. Crucially, learning trials were organized in different ways depending on the condition. In Exp. 1, the distractors from the two distributions were either mixed together (*Baseline*), colored differently (*Color*) or separated into different halves of the visual field (*Spatial*, see details below). On test trials, we randomly varied the similarity of the current target to non-targets from preceding trials (Figure 1B) with the aim of understanding how observers represent complex heterogeneous stimuli such as visual search distractors. The distractors on test trials were always from a Gaussian distribution centered at 60 to 120° relative to the current target. We assumed that during the learning trials observers encode the distractors and the distractor representation can be revealed by the response times on search trials. This would be consistent with our previous results where we have shown how response times follow the shape of the probability density function of the distractors, whether they are Gaussian or uniform, leftwards or rightwards skewed (Chetverikov et al., 2016, 2017d), or even bimodal versus uniform (Chetverikov et al., 2017c, 2020). To lay the groundwork for the analyses of empirical data, we first modeled the relationship between the distractor representations and response times in a Bayesian observer model.

### 2.1 | Bayesian observer model

How do behavioral responses depend on distractor representations from previous trials? To answer this question and to reconstruct distractor representations from the behavioral responses of our observers, we built a Bayesian memory-guided observer model linking observers' internal representations of distractors to response times (Figure 2A). This model is described here in short while the full description is available in the Methods.

We first describe a general structure of the model shown in Figure 2A. In the model, the observer had to locate a target among a set of distractors and indicate if it was in the top or the lower part of the stimuli matrix of size $N = 36$. The features (e.g., orientation) of each stimulus $s_i$ at locations $i = 1 \ldots N$ in the stimuli matrix are determined by the target location, $L_T$, and the parameters of the target feature distribution, $p\left(s_i \mid L_T = i\right)$ and of the distractor feature distribution, $p\left(s_i | L_T \neq i\right)$. For each trial, these parameters are used to generate stimulus $s_i$ at each location. At each moment in time $t$ within a trial, the observer obtains sensory samples or observations $x_{i,t}$ at each location. These
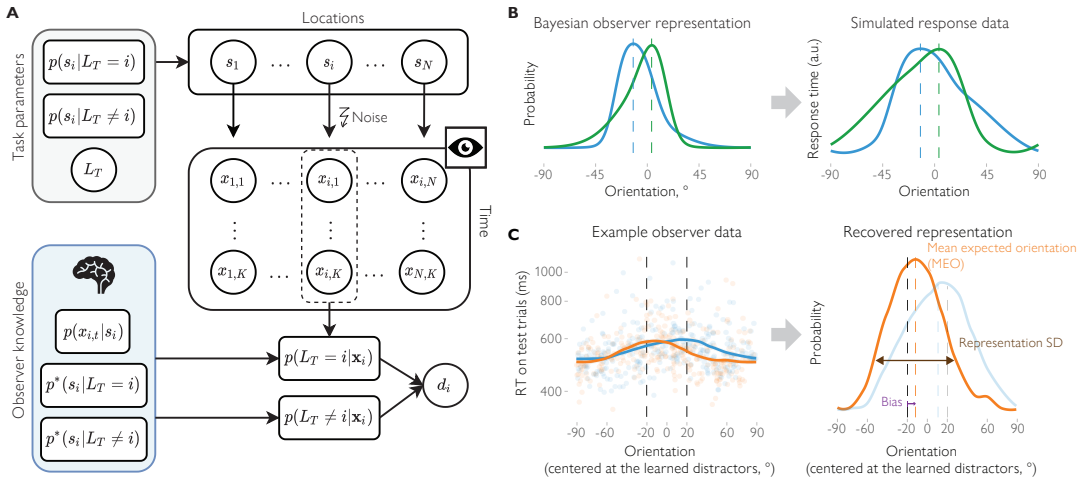
**FIGURE 2** The Bayesian observer model provides a way of reconstructing distractor representations. **A**: The Bayesian observer model. The stimuli $s_1 \ldots s_N$ at different locations are generated on each trial based on task parameters: the target feature distribution $p(s_i|L_T = i)$, the distractor feature distribution, $p(s_i|L_T \neq i)$, and the target location $L_T$. At each moment in time $t$ within a trial and for each location $i$, observers obtain samples of sensory observations $x_{i,t}$ corrupted by sensory noise, $p(x_{i,t} \mid s_i)$. Using knowledge about the sensory noise distribution and the approximation of feature distributions for targets and distractors, $p^*(s_i \mid L_T = i)$ and $p^*(s_i \mid L_T \neq i)$, observers compute probabilities that the sensory observations at a given location correspond to the target, $p(L_T = i \mid \mathbf{x}_i)$, or a distractor, $p(L_T \neq i \mid \mathbf{x}_i)$. 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 based on the information acquired from previous trials $p^*(s_i \mid L_T \neq i, \theta_{prev})$ (see text for more details; different example distributions are shown in blue and green). Crucially, there is a monotonic relationship between the two, with response times on test trials increasing as distractor probability increases. **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).

observations are not identical to the stimuli because of sensory noise that has a probability distribution $p(x_{i,t} \mid s_i)$. In other words, a given stimulus might result in different sensory responses, and, conversely, a given sensory observation might correspond to different stimuli.

Crucially, we do not assume that either the task parameters, such as $p(s_i|L_T = i)$ and $p(s_i|L_T \neq i)$, or the stimuli are known to the observer. However, the observer knows the parameters of sensory noise $p(x_{i,t} \mid s_i)$ and has an approximate knowledge of the target and distractor distributions denoted with asterisk, $p^*(s_i \mid L_T = i)$ and $p^*(s_i \mid L_T \neq i)$, that could be further separated into the knowledge based on the previous and on the current trial, e.g., $p^*(s_i \mid L_T \neq i, \theta_{prev})$ and $p^*(s_i \mid L_T \neq i, \theta_{curr})$ for distractors with $\theta_{prev}$ and $\theta_{curr}$ corresponding to the latent variables describing the parameters of the previous and the current trial, respectively. That is, in contrast to the traditional normative (ideal observer) models, our observer is not omniscient and does not know what was the distribution of distractors in the previous or the current trial. We assume instead that the observer has learned

some approximation of the distractor distribution from previous trials and combined it with the information about the current trials to improve search efficiency.

Using this knowledge, the observer aims to find the target by comparing for each location the probability that the sensory observations are caused by a target present at that location, $p\left(L_T = i \mid \mathbf{x}_i\right)$ against the probability that they are caused by a distractor, $p\left(L_T \neq i \mid \mathbf{x}_i\right)$:

$$d_i = \frac{p\left(L_T = i \mid \mathbf{x}_i\right)}{p\left(L_T \neq i \mid \mathbf{x}_i\right)} \tag{1}$$

where $\mathbf{x}_i = \{x_{i,1}, x_{i,2}, \ldots, x_{i,t=K}\}$ are the samples obtained for location $i$ up until a decision threshold is reached.

How can the observer estimate the probabilities in Eq. 1? Here, we would focus on the distractor-related part in the denominator but similar derivations can be done for the target. First, following the Bayes rule, the posterior probability that a distractor is at a given location is proportional to the likelihood of samples being drawn from the distractor distribution:

$$p\left(L_T \neq i \mid \mathbf{x}_i\right) \propto p\left(\mathbf{x}_i \mid L_T \neq i\right) \tag{2}$$

Assuming that the samples are independent in time their probability in log-space is equal to a sum of individual probabilities:

$$\log p\left(\mathbf{x}_i \mid L_T \neq i\right) = \sum_{t=1}^{K} \log p\left(x_{i,t} \mid L_T \neq i\right) \tag{3}$$

Importantly, the probability of a single observation corresponding to a distractor can be found by integrating over all possible stimuli values:

$$p\left(x_{i,t} \mid L_T \neq i\right) = \int p\left(x_{i,t} \mid s_i\right) p^*\left(s_i \mid L_T \neq i\right) ds_i \tag{4}$$

In other words, the observer combines the knowledge about sensory noise and the distractor distribution to estimate that a given sensory observation corresponds to distractors.

Notably, we are mainly interested in the test trials where the parameters of the current trial are independent of the parameters of the previous trials, hence:

$$p^*\left(s_i \mid L_T \neq i\right) \propto p^*\left(s_i \mid L_T \neq i, \boldsymbol{\theta}_{prev}\right) \tag{5}$$

To reiterate, 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 $p^*\left(s_i \mid L_T \neq i, \boldsymbol{\theta}_{prev}\right)$ matters for relative changes in response times.

Notably, if sensory observations are obtained with high frequency and sensory noise is low relative to the uncertainty in distractor representations, the log-sum of probabilities can be approximated as:

$$\sum_{t=1}^{K} \log p\left(x_{i,t} \mid L_T \neq i\right) \approx K\left(\log p^*\left(s_i \mid L_T \neq i\right) + C\right) \tag{6}$$

In words, if many samples are acquired at a given location, the log-probability that they are caused by a target is approximately equal to the number of samples obtained, times the log-probability that a stimulus at this location is a distractor.

Using Eq. 1, 5 and 6 (and analogous equations for target representations), the observer then can estimate the decision variable for all the samples obtained as:

$$\log d_i \propto K \left( \log p^* \left( s_i \mid L_T \neq i, \boldsymbol{\theta}_{prev} \right) - \log p^* \left( s_i \mid L_T \neq i, \boldsymbol{\theta}_{prev} \right) \right) \tag{7}$$

with constants subsumed under the proportionality sign. In words, the decision variable when a decision is made is proportional to a difference in the amount of evidence that a stimulus is a distractor and that it is a target times the number of samples obtained.

Finally, assuming that target and distractor representations are independent and noting that the response time is proportional to a number of observations needed to reach a decision, $RT \propto K$, the observer representation of distractor features learned from previous trials is related to response times:

$$RT \approx \frac{C_1}{C_0 - \log p^* \left( s_i \mid L_T \neq i, \boldsymbol{\theta}_{prev} \right)} \tag{8}$$

where $C_0$ and $C_1$ are constants (see details in Methods). In words, there is an inverse relationship between response times and the approximate likelihood that a given stimulus is a distractor, $p^* \left( s_i \mid L_T \neq i, \boldsymbol{\theta}_{prev} \right)$, with the information obtained from previous trials described by a set of latent parameters, $\boldsymbol{\theta}_{prev}$.

While this decision model is relatively simple, it provides a good intuition for observer behavior in the task (a more optimal model is provided in the Supplement but the conclusions do not depend on model choice). The model does not make any assumptions of how the observer learns these parameters. However, it shows that when the probability that a stimulus at a given location (e.g., a test target) is a distractor is lower, response times are lower as well, and vice versa.

This model provides an important insight, namely, that observers' representations are monotonically related to response times (Figure 2B). Hence, the monotonic relationship between the distribution parameters (mean, standard deviation, and skewness) reconstructed from RTs and from the true representation parameters would hold under any other monotonic transformation (for example, if RTs are log-transformed and the baseline is subtracted as we do in our analyses; see also Figure S2). In other words, the analysis above shows how response times can be used to approximately reconstruct observers' representations of distractors and estimate their parameters.

## 2.2 | Binding orientation probabilities to locations and colors

Having shown how observer response times can be related to the distractor representations, we now turn to the empirical data. Observers' response times to different test targets allow us to infer which orientations were the most difficult to find, resulting in the longest response times. As explained above, this methodology has enabled us to reveal how observers can represent distractor distributions in surprsising detail (Chetverikov et al., 2016, 2017d, 2020). Crucially, we can then reconstruct observers' representations of the probability distributions during learning trials (see Methods).

The experiments differed in the structure of the learning trials. There were three conditions in Experiment 1. The learning trials in the *Spatial* condition were organized so that distractor distributions in the left and the right hemifield differed to mimic the clustering of similar visual stimuli in the real world. In the *Color* condition, instead of spatial
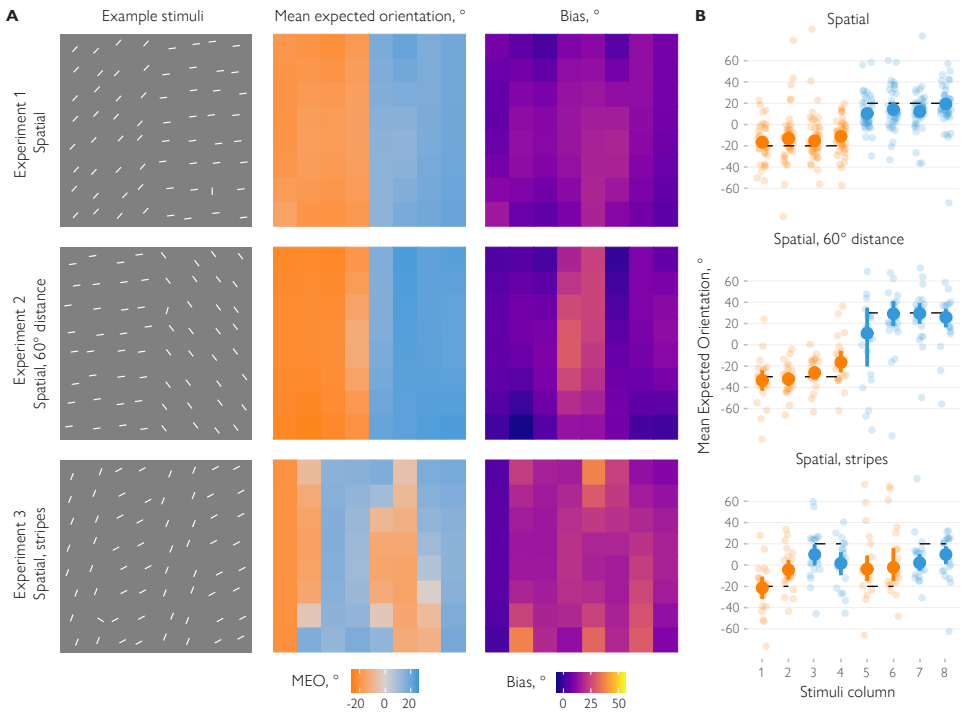
**FIGURE 3** Spatial structure of probabilistic representations. **A**: Example stimuli (left column), recovered mean expected orientations (middle column) and the across-distribution biases in mean expected orientations relative to the true orientations at a given location (right column). The stimuli 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.

grouping, different distractor subsets were grouped by color while individual items were randomly distributed. Finally, in the *Baseline* condition, items from the two distributions had the same color and were randomly distributed (Figure 1C).

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

We found that in the Spatial condition, observers' representations in each hemifield followed the actual physical

distractor distribution (Figure 3). The estimated MEO relative to the overall mean was *M* = -14.02° (*SD* = 6.02) and *M* = 14.90° (*SD* = 5.14) for probes for clockwise (CW) and counterclockwise-shifted (CCW) distributions, respectively. The difference in MEO between the two distributions was much greater than zero (*b* = 28.94°, 95% HPDI = [25.34, 32.56], *BF* = $6.35 \times 10^{17}$) showing that observers expected different orientations in different hemifields. We then computed the across-distribution bias by recoding the errors in MEO relative to the true mean for each distribution so that positive values correspond to shifts towards the other distribution. That is, the bias here represents by how much observers' expectations deviated from the true mean orientation at a given location towards the mean orientation at the other location. For both hemifields there was a significant bias towards the other hemifield (*M* = 5.52°, 95% CI = [1.86, 9.14]). This shows that while observers represent the spatial separation between the two distributions, signals from the other hemifield still influence their responses.

But does spatial separation help observers track the feature probabilities? In the Baseline conditions, locations of the CW and CCW distributions were chosen randomly for each learning trial. We repeated the analysis described above, comparing the response to test targets at the location that had CW and CCW orientations on the immediately preceding trials. We expected to find stronger across-distribution biases as there was no separation between the distributions across trials. Importantly, the across-distribution bias was larger in the Baseline (bias *M* = 11.35°, 95% CI = [7.71, 15.00]; Figures 3 and 4) than the Spatial condition (effect of condition *M* = 5.84, 95% CI = [1.10, 10.58], *BF* = 108.24). In other words, the representations for each distribution were closer to the overall distribution mean in the Baseline than the Spatial condition. This argues that when the learned distributions are consistently presented at separate locations, observers can track them better than when they are randomly distributed.

Do observers integrate information about orientation probabilities and color? In the Color condition, the locations of the test targets were counterbalanced with respect to their colors, so we should only find differences in MEO if observers formed an association between color and orientation. Indeed, we found that the MEOs for the two distributions differed (*b* = 7.35, 95% HPDI = [1.30, 13.06], *BF* = 148.04) although across-distribution biases were stronger (*M* = 16.30, 95% CI = [12.66, 19.86]) than in the Spatial condition (the difference between conditions *M* = 10.78, 95% CI = [5.99, 15.54], *BF* = $6.56 \times 10^4$). This means that if observers saw yellow lines shifted CW and blue lines shifted CCW relative to the overall distractor mean during learning trials, they learned this association which affected their response times on subsequent test trials. Importantly, this demonstrates that observers can integrate information about likely orientations with information about other features (in this case color), even if there is no spatial information to guide this integration.

## 2.3 | Encoding orientation probabilities at different spatial scales

Having established that observers associate information about most likely orientations with specific locations or colors, we then asked if we can uncover the origins of the observed biases by assessing the recovered representations in the Spatial condition in more detail (for this and later analyses, we increased the sensitivity of our analyses by combining the data from the Spatial group in Experiment 1 with an additional participant group that performed the same task; see Methods). We computed MEO using the aggregated data from all participants for each location in the stimuli matrix in this condition. As Figure 3 shows, across-distribution biases were stronger closer to the boundary between the two hemifields. We then tested this observation by directly comparing MEOs for test trials with targets presented at the boundary (two central columns) between the hemifields against other test trials. We found that the bias was significantly larger at the boundary between the two distributions than in the other columns (*M* = 4.80° (*SD* = 6.99) and *M* = 9.04° (*SD* = 11.36), *b* = 4.23, 95% HPDI = [0.21, 8.32], *BF* = 42.34; Figure 3B). However, the biases were also significantly above zero outside the boundary (*BF* = 248). This suggests that the distribution representations are not

homogeneous and influence each other strongly when they are close in space, but this mutual influence also extends outside the immediate neighboring locations (see Discussion).

## 2.4 | Bias strength depends on similarity and spatial arrangement

In two follow-up studies, we further investigated observers' representations of spatially-grouped heterogeneous stimuli. In Experiment 2, we tested whether the similarity between the distributions along the tested feature dimension (orientation) affects the strength of the across-distribution biases. Recent studies suggest that similarity is an important factor determining whether the information is pooled or not at different levels of perceptual processing (e.g., Coen-Cagli et al., 2015; Herrera-Esposito et al., 2021; Manassi et al., 2012; Qiu et al., 2013; Utochkin et al., 2018). We hypothesized that the bias should be stronger when the stimuli from the two distributions are more likely to have the same cause in the external world. For example, the boundary effect in Experiment 1 might occur because stimuli that are close in space are more likely to belong to the same object. By the same reasoning, if the two distributions are less similar, they are less likely to have the same cause, and the biases should be weaker.

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

In Experiment 3, we tested whether an even more complex spatial arrangement would allow us to recover the "map" of observers' expected orientations. To this end, the stimuli were organized in "stripes" of two matrix columns with two different distributions from Experiment 1 (with means separated by 40°) positioned at odd and even stripes (counterbalanced across blocks, Figure 3A). We found that observers expected clockwise-rotated orientations ($M$ = 6.20°, $SD$ = 9.91) at locations of stripes rotated 20° clockwise relative to the overall mean and counterclockwise-rotated orientations ($M$ = -11.04°, $SD$ = 17.11) at other stripe locations. However, the across-distribution bias ($M$ = 11.70°, $SD$ = 7.52) was stronger than in the Spatial condition in Experiment 1 ($b$ = 5.90, 95% HPDI = [2.50, 9.33], $BF$ = 4.30). This demonstrates that while separating distributions in space helps observers track distributions (as shown in Experiments 1 and 2), the effects of spatial organization decrease as the organization becomes more complex.

## 2.5 | Higher-order parameters of probabilistic representations

Next, we asked whether observers' representations contain more information about the distributions than just their average? We used the reconstructed distractor representations (Figure 4A) to estimate their circular standard deviation and circular skewness. The former corresponds to the expected variability among distractors, while the latter quantifies their symmetry.

First, we hypothesized that if the variability of the distributions is encoded, then the expected variability would be higher when the distractor distributions are less well separated. Indeed, we found that observers' expectations about distractor variability differ between conditions ($BF = 2.03 \times 10^5$) with lower SD when the distractors were separated by hemifields ($M$ = 33.3, 95% HPDI = [32.2, 34.4] for the Spatial condition with 40° separation and $M$ = 32.7, 95% HPDI = [31.1, 34.2] for 60° separation) compared to other conditions ($M$ = 35.9, 95% HPDI = [34.4, 37.5] in the color condition, $M$ = 34.4, 95% HPDI = [32.9, 35.9] for the stripes arrangement condition). When the two distributions
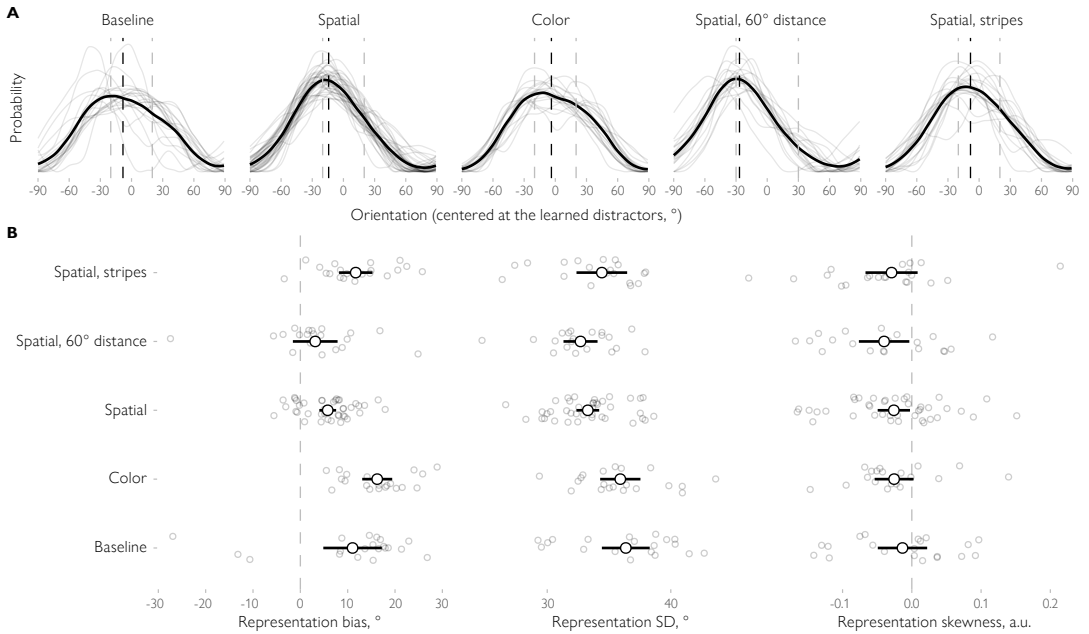
**FIGURE 4** Recovered average representations and their parameters across experiments and conditions. **A**: The black curves show the average representation while representations for individual observers are shown in light gray. Dashed vertical lines show the mean of the representation (black) and 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° or -30°) while the other mean is counterclockwise (20° or 30° 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.

were less well separated, observers were more uncertain in their estimates, leading to distractor representations with higher SD's (Figure 4B).

We also expected that the distribution presented at the tested location or in the tested color would weigh more highly in the resulting representation, causing an asymmetry. Alternatively, if observers only use the mean and variance to encode the distribution (as assumed by "summary statistics" accounts), then the represented distribution should be symmetric. We found that observers' representations were asymmetric in all conditions, with a higher probability mass at the side corresponding to the distribution presented at the tested location or in the tested color, $M$ = -0.03, 95% CI = [-0.04, -0.02]. Notably, however, no differences between conditions were found, $BF = 1.99 \times 10^{-6}$, indicating that symmetry is not affected by the way the distributions are organized in the display. In sum, observers represent not only the average stimulus values but also their variability, and the representations are skewed towards distributions presented at other locations or in different colors.

# 3 | DISCUSSION

Our main hypothesis was that observers extract information about probabilities of visual features from heterogeneous stimuli and bind the resulting probabilistic representations with locations on the one hand and other features on the other. Our results support both these proposals very clearly. Importantly, this demonstrates for the first time, how the visual system can build probabilistic representations of the visual world by extracting information about the features of complex heterogeneous stimuli.

A visual search task allowed us to uncover representations of heterogeneous distractors. We formulated a Bayesian observer model and demonstrated analytically and through simulations that response times are a monotonic function of observers' expectations about distractor orientations, supporting earlier empirical findings (Chetverikov et al., 2016, 2017b, 2020; Chetverikov et al., 2019; Hansmann-Roth et al., 2019; Hansmann-Roth et al., 2021; Tanrıkulu et al., 2020, 2021). Using this knowledge, we were able to estimate the characteristics of observer representations – their means, precision, and skewness – and to assess how they vary depending on whether observers can associate them with locations or with other, task-irrelevant features, such as color.

We found that observers both encode and combine feature distributions in scenes containing two different distributions. The representations generally follow the physical distribution of the stimuli for a given location or a given color, but importantly, observers are also biased towards the other distribution. The strength of the bias depends on the degree of separation between the distributions. When the distributions were separated in space, observers' representations of one distribution were less influenced by the other distribution, compared to when they were separated by color or were intermixed (Baseline condition). Furthermore, as we found in Experiment 3, more complex spatial arrangements ("stripes") increased the biases towards the other distribution. In sum, observers bind probabilistic representations of visual features to locations and other features, but such binding is not impenetrable, reminiscent of 'illusory conjunctions' of discrete feature values (Treisman & Schmidt, 1982).

We were then able to recover the representation of the distribution at different spatial scales. We found that for spatial separation, the biases are stronger at the boundary between the two distributions. This is reminiscent of the hierarchical organization of information about feature probabilities within a scene proposed for perceptual ensembles (Alvarez, 2011; Haberman & Whitney, 2012). Such hierarchical ensemble models suggest that observers represent information about feature probabilities at different levels: for example, the orientation statistics at a particular location are combined to form a representation for a group of items, which are, in turn, combined to form an overall ensemble representation. Our results agree with this idea: the stimuli observers expect at a given location depend not only on what was previously shown at this location but also on stimuli presented at other locations. Crucially, biases were also present for the Color condition as well as for the non-boundary locations in the Spatial condition of Experiment 1. This indicates that the results cannot be explained by purely local summation of the inputs. It remains to be tested whether there are actual separable representations of probability distributions at different levels, or just a unified spatio-featural map guiding observer responses.

We hypothesized that the representations should be more biased by each other when they are more likely to have the same cause in the external world. This could provide a normative explanation for the boundary effect: sensory input from adjacent locations is likely to be caused by the same object and should therefore be integrated while locations far away from each other should be treated separately. Similarly, for example, in multisensory integration studies, auditory and visual signals are less likely to be integrated when there is a large discrepancy in their locations (Körding et al., 2007; Shams & Beierholm, 2010). However, in Experiment 1, we found across-distribution biases at locations far from the other distribution. We reasoned that this is because the stimuli themselves are similar enough to be potentially caused by the same object, and the inputs are therefore integrated even from non-neighboring locations. In

Experiment 2, we tested this explanation by asking if the similarity between the distributions themselves in the tested feature domain (orientation) also plays a role. We found that when the distributions were made more dissimilar, the biases were observed only at the boundary between the distributions but not at other locations. That is, observers no longer take into account the input from non-neighboring locations, when stimuli are dissimilar. Speculatively, introducing longer learning streaks could also help to reduce the bias by increasing precision of the representations (Chetverikov et al., 2017c). This supports the proposed normative explanation and suggests that the principles of information integration for heterogeneous visual inputs are the same as for other cases, such as multisensory integration or estimation of complex visual features (Landy et al., 2011).

We then tested if observers represent more than just the mean distractor orientation. We found that observers represent the distractor variability (i.e., the standard deviation or width of their representations), which varies in a predictable fashion with the separability between distractor distributions. When distractor distributions are poorly separated (e.g., only by color or are organized in 'stripes'), their representations are wider, indicating more uncertainty. Furthermore, the representations are asymmetric where the tail of the distribution corresponding to the orientations matching the tested location or color is fatter. While we are agnostic to the specific mechanisms of how the information is integrated across different parts of the visual field, speculatively, such asymmetric distributions can be seen as an output of a hypothetical weighting process. As a computational abstraction, the resulting representation can be seen as a normalized sum of basis functions (similar to a kernel density estimator). The weight of a certain basis function could depend on how well it matches the stimuli across the visual field (i.e., how many distractors had a certain orientation) and their relevance to the current goals. The presence of skewness indicates that observers do not simply represent the distractors with a (biased) mean and variance, their representations have a complex shape with more relevant information (e.g., previous orientations at a tested location) weighted higher and less relevant information (e.g., previous orientations at the other locations) having lower weight, but still influencing the outcome. However, we do not see the effect of condition on the distribution asymmetry, which would be expected if the matching and nonmatching parts of the distribution were combined as a weighted mixture with weights depending on the degree of separation. The absence of the condition effect might be related to the greater difficulty of precisely estimating the amount of skewness as opposed to mean or variance. Nevertheless, the overall skew in the representations is indicative of how sophisticated the learning can be where various factors, such as the amount of information about the underlying probability distribution and task-relevance of the stimuli in each case, are taken into account in the representation, and how this determines how different parts of the display are weighted.

These findings indicate that observers represent information about distractor features as a probability distribution rather than only in terms of the summary statistics, in contrast to popular ideas of simple "summary statistics". For example, Treisman (2006) argued that statistical processing is a distinct mode of perceptual and attentional analysis of stimulus sets. She proposed that because of limited attentional capacity statistical summaries are generated that include the mean, variance, and perhaps the range. These summaries enable rapid assessment of the general properties and layout of natural scenes (Chong & Treisman, 2005; Emmanouil & Treisman, 2008). Similarly, Rahnev (Rahnev, 2017; Yeon & Rahnev, 2020) argued that observers represent only a summary consisting of the most likely stimulus and the associated strength of evidence, and Cohen et al. (2016) used summary statistics to explain the richness of consciousness experience. Our results argue against such views, since the representations that are bound together are far more detailed than this implies. That is, the brain might instead approximate the visual input by using a complex set of parameters to provide accurate descriptions of feature probabilities (Freeman & Simoncelli, 2011; Rosenholtz, 2020).

A recent finding may provide insights into why summary statistic accounts have been so poular. Hansmann-Roth et al. (2021) reasoned that optimal behavior requires the encoding of full feature distributions, not only summaries,

but observers might be unable to explicitly report the full distribution. This is analogous to how difficult it might be to verbally describe the variety of colors of an apple without resorting to simplifications (see Figure 1A). Hansmann-Roth et al. tested observers' representations both implicitly and explicitly and while explicit judgments were limited to the mean and variance of feature distributions, implicit measures revealed detailed representations of the same distributions. More information was therefore available to observers than studies of summary statistics, that have mostly relied on explicit measures, have indicated. Crucially, Hansmann-Roth et al. were able to uncover why this is: revealing these detailed representations requires implicit methods where representations are probed by assessing the effects of role reversals between targets and distractors as we do here.

Can the results observed in our study be explained by simple sensory adaptation mechanisms? Sensory adaptation is a well-studied phenomenon: at a neural level, exposure to a stimulus alters neural responses to subsequently presented stimuli, while behaviorally adaptation often results in a repulsive bias with estimates of, for example, orientation in the adjustment task shifted away from the adapter (Clifford et al., 2007; Schwartz et al., 2007). In our task, observers were exposed to a certain distractor distribution within each miniblock, so their behavior could be influenced by adaptation. We tested this idea by developing a variation of a model reported here assuming that the observer does not encode the knowledge about the distractors directly ($p^*\left(s_i \mid L_T \neq i, \theta_{prev}\right)$ is flat) but instead the sensory space is warped to efficiently encode the distractor distribution (Stocker & Simoncelli, 2005; Wei & Stocker, 2012, 2015). The results (not shown) indicated that the search would be the most efficient when a target matches previous distractors, contrary to our findings and the well-known role-reversal effects in visual search literature (Kristjánsson & Driver, 2008). The intuition here is that when targets and distractors are on the opposite sides of the adapter, they are actually pushed together due to the circularity of orientation space. Furthermore, recent studies using a similar behavioral task by Rafiei, Chetverikov, et al., 2021; Rafiei, Hansmann-Roth, et al., 2021 and Pascucci et al., 2022 also speak against the involvement of adaptation. Rafiei et al. looked at how the perception of a current search target (2021) and a neutral line (2021) is affected by distractors and previous targets. Importantly for the current issue, they found that when distractors are similar to the test item, its estimates are pulled towards the distractors, while the adaptation profile is generally repulsive for similar items. In a recent study by Pascucci, Ceylan, and Kristjánsson (2022), the authors found that when observers passively viewed an array of lines that all came from the same distribution (no singleton), no learning of the distribution occurred. In summary, both the modelling results and the empirical data suggest that sensory adaptation cannot explain our findings.

Where does the uncertainty in the distractor representation come from and how the observers learn about it? Our study is agnostic on this issue. We do not assume any parametric form for this representation (that is, we do not assume that observers represent, for example, the mean or variance). Note, however, that it is unlikely that a single-trial representation is a simple Gaussian and only when aggregating across trials the skewness appears. This idea was tested in a recent study by Chetverikov et al. (2020), where observers have to find two targets on each trial with a mixture distribution similar to the Baseline condition used here. By comparing a model assuming that the distribution is approximated as a single-peaked and a model assuming that the distribution is approximated as a mixture of two parts, we found that the hypothesis of a simple Gaussian representation is not supported by the data. It might well be, however, that at a more detailed level, for example, at a level of a single location or a single moment in time, the representation is a simple one. Yet, from the computational perspective, these simpler representations have to be combined into an aggregated representation of distractors shown here and in our previous studies to be corresponding to a probability distribution of distractors with more details that the simple summary statistics would allow.

In our experiments, observers learn the distractor feature by combining inputs from heterogeneous stimuli across several trials in each block, and it can be argued that this is different from perceiving a single stimulus on a single trial.

However, the visual cortex aggregates information on many different timescales (de Lange et al., 2018). Even on a single trial, perception unfolds in time and at each moment is dependent on what has been seen before. And even for a simple stimulus, the visual cortex receives inputs from many retinal neurons that are affected by processing noise, potentially indistinguishable from the input from varying features. Indeed, this is why stimulus variability ('external noise') is often used to manipulate visual uncertainty (Barthelmé & Mamassian, 2009; Hénaff et al., 2020). We therefore believe that distinguishing "simple" and "complex" perception is impossible. However, our results clearly show that information about feature probabilities is available for visually-guided behavior.

## 4 | SUMMARY

Taken together, our results show that observers can not only encode probabilities of features from heterogeneous stimuli in detail, but also integrate them with both locations and other features that have different distributions. These results arguably represent the strongest support yet for the long-standing idea that the brain builds probabilistic models of the world (Chetverikov et al., 2017a; Fiser et al., 2010; Knill & Pouget, 2004; Orhan & Ma, 2015; Rao et al., 2002; Sahani & Dayan, 2003; Tanrıkulu, Chetverikov, Hansmann-Roth, et al., 2021) and show that probabilistic representations can serve as building blocks for object and scene processing. Notably, such representations are not simply limited to summary statistics (e.g., a combination of mean and variance; Cohen et al., 2016). Our results also indicate that observers do not represent physical stimuli precisely, but instead construct an approximation influenced by input from other stimuli. This probabilistic perspective stands in sharp contrast to views where discrete features of individual stimuli are *either* bound together to form objects or processed "statistically" (Rosenholtz, 2020; Treisman, 2006). Instead, we suggest that the probabilistic representations are automatically bound to locations and other features since such binding occurred even though it was not required in the task. Probabilistic representations are therefore not acquired in isolation but constitute an integral part of perception.

## 5 | METHODS

### 5.1 | Participants

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

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

## 5.2 | Procedure

In *Experiment 1*, each participant performed a search task in five conditions. In each condition on each trial, observers were presented with 8×8 matrices of 64 lines (line length: 0.71° of visual angle; matrix size: 16×16°; uniform noise of ±0.5° was added to each line coordinate). The goal was to find the odd-one-out line whose orientation differed most from the others. Sessions were divided into miniblocks of 5 to 7 learning trials followed by 1 or 2 test trials (the number of trials chosen randomly for each block; the variation in the number of trials was introduced to decrease the effect of temporal expectations, Shurygina et al., 2019. During learning trials, the overall mean of distracting items stayed the same within each miniblock (but varied randomly between miniblocks) with half of the distractors drawn from one distribution and the other half from another distribution with the properties of distributions differing between conditions:

*Baseline*: two truncated Gaussian distributions with SD = 10° and range of 40°, with means separated by 40° (±20° relative to the overall mean), all stimuli had the same color (white); half of the distractors were drawn randomly from one distribution, half from another, then they were positioned randomly within the stimuli matrix and then a randomly chosen distractor was replaced with a target.

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

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

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

On test trials, distractors came from a single Gaussian distribution with SD = 10° (range-restricted in the same way as described above) and a mean sampled from a 0-180° uniform distribution, while target orientation was determined in the same way as on the prime trials (that is, selected randomly from 60 to 120° range relative to the current distractor mean). In the color condition, half of the lines from that distribution were blue, half were yellow.

The Baseline condition had 2304 trials, while the Spatial and Color conditions had 5376 trials each with the higher number of trials used in the latter case to counterbalance additional factors (distribution type combinations). The trials were split into two (for Baseline) or four sessions (other conditions) with a break for rest halfway within each session. Observers participated in each session at a separate time depending on availability but with a break of no less than two hours between sessions and no more than two sessions within a day. The order of sessions with respect to conditions was counterbalanced with a Latin square design.

*Experiments 2 and 3* followed the same general procedure as the Spatial condition of Experiment 1. In Experiment 2 the means of the distributions were separated by 60° (±30° relative to the overall mean) instead of 40° in Experiment 1. In Experiment 3, the two distributions were separated by 40°, as in Experiment 1, but arranged in "stripes" so that the lines drawn from the first distribution were positioned in the 1st, 2nd, 5th, and 6th columns of the stimuli matrix while the other columns were populated with lines from the second distribution. In both experiments, each participant took part in two sessions of 1536 trials each with a rest period halfway within each session.

## 5.3 | Data processing

For our main analyses of interest, incorrect responses were excluded, and response times were log-transformed and centered by subtracting the mean for each participant. Then, to reduce the noise in RT measurements, spatial and featural confounders were removed (the results remain similar when no corrections are applied). First, the effect of the distance between target locations on consecutive trials and the effect of the target location were removed by regressing out the fifth-degree polynomials of the absolute distance (in degrees of visual angle) between the target locations on the current and the previous trials and the current targets horizontal and vertical coordinates. Then, we also removed potential influences from the well-known oblique effect (the search speed differs between oblique and cardinal stimuli Chetverikov et al., 2017a; Wolfe et al., 1999 by regressing out the fifth-degree polynomials of target and distractor obliqueness computed as an absolute distance in degrees to the nearest cardinal orientation. The regression was run separately for each experiment and condition.

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

To convert the noisy response times into probabilities, we first smoothed RT as a function of the test target and previous distractor mean using the local regression approach (a generalization of the moving average) for each observer in each condition. To account for circularity, we appended 1/6 of the data from each end of the orientation space to the opposite end before smoothing. In analyses applied to each stimulus location, we further assumed that RTs are a smooth function of the stimuli matrix row within the local regression while columns of the stimuli matrix were treated independently. We then transformed a smoothed RT function into a probability mass function by subtracting the baseline and normalizing to one. Finally, we computed the parameters of the recovered probabilistic representation: the mean expected orientation (circular mean), circular standard deviation, and circular skewness as defined by Pewsey (Pewsey, 2004). Note that under the hypothesized Bayesian observer model, the estimated standard deviation and skewness are monotonically related to the true parameters of the distractor representation but are not identical to it (additionally confirmed in simulations, Figure S2).

Unless stated otherwise, we used Bayesian hierarchical regression with the *brms* (Bürkner, 2017) package in R. Note that while we include Bayes factor values in the description of the results, we were mostly interested in measuring the effects of the variables of interest in our models, and hence the models included the default flat (uniform) priors for regression coefficients. Given that Bayes factors are prior-dependent, we believe that the information provided by the 95% highest-density posterior intervals (HDPI) is more useful for judging the results than the Bayes factors. To make sure that the conclusions are not dependent on the particular analytic approach, we repeated the analyses using the conventional frequentist statistical test with the same results (the report using this approach is provided alongside the data in an online repository, see data availability statement).

## 5.4 | Bayesian observer model

In our experiments, participants located a target among a set of distractors and indicated if it was in the upper or the lower part of the stimuli matrix. On each trial, the experimenter sets the task parameters, namely, parameters of the target distribution, $p(s_i|L_T = i)$, and parameters of the distractor distribution, $p(s_i|L_T \neq i)$, for each location $i = 1 \ldots N$ in the stimuli matrix as well as the target location, $L_T$. These parameters were then used to generate the stimuli $s_i$ at each location.

Neither the task parameters nor the stimuli are known to the Bayesian observer. Instead, at each moment in time

$t$ within a trial, the observer obtains sensory observations at each location, $x_{i,t}$. These observations are not identical to the stimuli because of the presence of sensory noise, $p\left(x_{i,t} \mid s_i\right)$. That is, a given stimulus might result in different sensory responses, and, conversely, a given sensory observation might correspond to different stimuli. We assume that the observations are distributed independently at each location and at each moment in time.

To make an optimal decision in a particular task, the observer needs to know the relationship between the sensory observations and the task-relevant quantities. For the visual search task used in our study, we assumed that observers compare for each location the probability that the sensory observations are caused by a target present at that location, $p\left(L_T = i \mid \mathbf{x}_i\right)$ where $\mathbf{x}_i = \{x_{i,1}, x_{i,2}, \ldots, x_{i,t=K}\}$ are the samples obtained for location $i$ up until the time $K$, against the probability that they are caused by a distractor, $p\left(L_T \neq i \mid \mathbf{x}_i\right)$:

$$d_i = \frac{p\left(L_T = i \mid \mathbf{x}_i\right)}{p\left(L_T \neq i \mid \mathbf{x}_i\right)} \tag{9}$$

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

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

$$p\left(L_T = i \mid \mathbf{x}_i\right) = \frac{p\left(\mathbf{x}_i \mid L_T = i\right) p\left(L_T = i\right)}{p\left(\mathbf{x}_i\right)} \tag{10}$$

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

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

$$\log d_i = \log\left(\frac{1}{N-1}\right) + \sum_{t=1}^{K} \log\left(p\left(x_{i,t} \mid L_T = i\right)\right) - \sum_{t=1}^{K} \log p\left(x_{i,t} \mid L_T \neq i\right) \tag{11}$$

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

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

$$p\left(x_{i,t} \mid L_T \neq i\right) = \int p\left(x_{i,t} \mid s_i\right) p^*\left(s_i \mid L_T \neq i\right) ds_i \tag{12}$$

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

$$p^*\left(s_i \mid L_T \neq i\right) \equiv p^*\left(s_i \mid L_T \neq i, \theta\right) = p^*\left(s_i \mid L_T \neq i, \theta_{prev}\right) p^*\left(s_i \mid L_T \neq i, \theta_{curr}\right) \tag{13}$$

with $\theta = \{\theta_{prev}, \theta_{curr}\}$ 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 $p^*\left(s_i \mid L_T \neq i, \theta_{prev}\right)$ 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 proportional 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 $p^*\left(s_i \mid L_T \neq i, \theta_{prev}\right)$.

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

$$\sum_{t=1}^{K} \log p\left(x_{i,t} \mid L_T \neq i\right) \approx K\left(E\left[\log p\left(x_{i,t} \mid L_T \neq i\right)\right]\right) \tag{14}$$

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

$$E\left[\log p\left(x_{i,t} \mid L_T \neq i\right)\right] \propto \log p^*\left(s_i \mid L_T \neq i\right) + C \tag{15}$$

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

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

$$K = \frac{\log B - \log\left(\frac{1}{N-1}\right)}{E\left[\log p\left(x_{i,t} \mid L_T = i\right)\right] - E\left[\log p\left(x_{i,t} \mid L_T \neq i\right)\right]} \tag{16}$$

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

$$RT \approx \frac{C_1}{C_0 - \log p^*\left(s_i \mid L_T \neq i, \theta_{prev}\right)} \tag{17}$$

and

$$\log p^*\left(s_i \mid L_T \neq i, \theta_{prev}\right) = C_0 - C_1 \frac{1}{RT} \tag{18}$$

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.

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

## Acknowledgements

## Data availability

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

## Conflict of interest

The authors declare no conflict of interest.

## References

Alvarez, G. A. (2011). Representing multiple objects as an ensemble enhances visual cognition. *Trends in cognitive sciences*, *15*(3), 122–31. https://doi.org/10.1016/j.tics.2011.01.003

Ariely, D. (2001). Seeing sets: Representation by statistical properties. *Psychological Science*, *12*(2), 157–162. https://doi.org/10.1111/1467-9280.00327

Attarha, M., & Moore, C. M. (2015a). The perceptual processing capacity of summary statistics between and within feature dimensions. *Journal of Vision*, *15*(4), 9. https://doi.org/10.1167/15.4.9

Attarha, M., & Moore, C. M. (2015b). The capacity limitations of orientation summary statistics. *Attention, Perception, & Psychophysics*, *77*(4), 1116–1131. https://doi.org/10.3758/s13414-015-0870-0

Attarha, M., Moore, C. M., & Vecera, S. P. (2014). Summary statistics of size: Fixed processing capacity for multiple ensembles but unlimited processing capacity for single ensembles. *Journal of Experimental Psychology: Human Perception and Performance*, *40*(4), 1440–9. https://doi.org/10.1037/a0036206

Balas, B. J., Nakano, L., & Rosenholtz, R. (2009). A summary-statistic representation in peripheral vision explains visual crowding. *Journal of vision*, *9*(12), 13.1–18. https://doi.org/10.1167/9.12.13

Barthelmé, S., & Mamassian, P. (2009). Evaluation of objective uncertainty in the visual system (K. Kording, Ed.). *PLoS Computational Biology*, *5*(9), e1000504. https://doi.org/10.1371/journal.pcbi.1000504

Block, N. (2018). If perception is probabilistic, why does it not seem probabilistic? *Philosophical Transactions of the Royal Society B: Biological Sciences*, *373*(1755). https://doi.org/10.1098/rstb.2017.0341

Bürkner, P.-C. (2017). Brms : An r package for bayesian multilevel models using stan. *Journal of Statistical Software*, *80*(1). https://doi.org/10.18637/jss.v080.i01

Chetverikov, A., Campana, G., & Kristjánsson, Á. (2016). Building ensemble representations: How the shape of preceding distractor distributions affects visual search. *Cognition*, *153*, 196–210. https://doi.org/10.1016/j.cognition.2016.04.018

Chetverikov, A., Campana, G., & Kristjánsson, Á. (2017a). Learning features in a complex and changing environment: A distribution-based framework for visual attention and vision in general. *Progress in brain research* (pp. 97–120). Elsevier. https://doi.org/10.1016/bs.pbr.2017.07.001

Chetverikov, A., Campana, G., & Kristjánsson, Á. (2017b). Set size manipulations reveal the boundary conditions of distractor distribution learning. *Vision Research*, *140*(November), 144–156. https://doi.org/10.1016/j.visres.2017.08.003

Chetverikov, A., Campana, G., & Kristjánsson, Á. (2017c). Rapid learning of visual ensembles. *Journal of Vision*, *17*(21), 1–15. https://doi.org/10.1167/17.2.21

Chetverikov, A., Campana, G., & Kristjánsson, Á. (2017d). Representing color ensembles. *Psychological Science*, *28*(10), 1–8. https://doi.org/10.1177/0956797617713787

Chetverikov, A., Campana, G., & Kristjánsson, Á. (2020). Probabilistic rejection templates in visual working memory. *Cognition*, *196*, 104075. https://doi.org/10.1016/j.cognition.2019.104075

Chetverikov, A., Hansmann-Roth, S., Tanrikulu, Ö. D., & Kristjánsson, Á. (2019). Feature distribution learning (FDL): A new method for studying visual ensembles perception with priming of attention shifts. *Neuromethods* (pp. 1–21). Springer. https://doi.org/10.1007/7657_2019_20

Chong, S. C., & Treisman, A. (2005). Statistical processing: Computing the average size in perceptual groups. *Vision research*, *45*(7), 891–900. https://doi.org/10.1016/j.visres.2004.10.004

Clifford, C. W. G., Webster, M. A., Stanley, G. B., Stocker, A. A., Kohn, A., Sharpee, T. O., & Schwartz, O. (2007). Visual adaptation: Neural, psychological and computational aspects. *Vision Research*, *47*(25), 3125–3131. https://doi.org/10.1016/j.visres.2007.08.023

Coen-Cagli, R., Kohn, A., & Schwartz, O. (2015). Flexible gating of contextual influences in natural vision. *Nature Neuroscience*, *18*(11), 1648–1655. https://doi.org/10.1038/nn.4128

Cohen, M. A., Dennett, D. C., & Kanwisher, N. (2016). What is the bandwidth of perceptual experience? *Trends in Cognitive Sciences*, *20*(5), 324–335. https://doi.org/10.1016/j.tics.2016.03.006

de Lange, F. P., Heilbron, M., & Kok, P. (2018). How do expectations shape perception? *Trends in Cognitive Sciences*, *22*(9), 764–779. https://doi.org/10.1016/j.tics.2018.06.002

Emmanouil, T. A., & Treisman, A. (2008). Dividing attention across feature dimensions in statistical processing of perceptual groups. *Perception & psychophysics*, *70*(6), 946–954. https://doi.org/10.3758/PP.70.6.946

Fiser, J., Berkes, P., Orbán, G., & Lengyel, M. (2010). Statistically optimal perception and learning: From behavior to neural representations. *Trends in cognitive sciences*, 119–130. https://doi.org/10.1016/j.tics.2010.01.003

Freeman, J., & Simoncelli, E. P. (2011). Metamers of the ventral stream. *Nature neuroscience*, *14*(9), 1195–1201. https://doi.org/10.1038/nn.2889

Girshick, A. R., Landy, M. S., & Simoncelli, E. P. (2011). Cardinal rules: Visual orientation perception reflects knowledge of environmental statistics. *Nature Neuroscience*, *14*(7), 926–932. https://doi.org/10.1038/nn.2831

Haberman, J., & Whitney, D. (2012). Ensemble perception: Summarizing the scene and broadening the limits of visual processing. In J. M. Wolfe & L. Robertson (Eds.), *From perception to consciousness: Searching with anne treisman* (pp. 339–349). Oxford University Press. https://doi.org/10.1093/acprof:osobl/9780199734337.003.0030

Hansmann-Roth, S., Chetverikov, A., & Kristjánsson, Á. (2019). Representing color and orientation ensembles: Can observers learn multiple feature distributions? *Journal of Vision*, *19*(9), 1–17. https://doi.org/10.1167/19.9.2

Hansmann-Roth, S., Kristjánsson, Á., Whitney, D., & Chetverikov, A. (2021). Dissociating implicit and explicit ensemble representations reveals the limits of visual perception and the richness of behavior. *Scientific Reports*, 1–12. https://doi.org/10.1038/s41598-021-83358-y

Hénaff, O. J., Boundy-Singer, Z. M., Meding, K., Ziemba, C. M., & Goris, R. L. (2020). Representation of visual uncertainty through neural gain variability. *Nature Communications*, *11*(1), 1–12. https://doi.org/10.1038/s41467-020-15533-0

Herrera-Esposito, D., Coen-Cagli, R., & Gomez-Sena, L. (2021). Flexible contextual modulation of naturalistic texture perception in peripheral vision. *Journal of Vision*, *21*(1), 1. https://doi.org/10.1167/jov.21.1.1

Knill, D. C., & Pouget, A. (2004). The bayesian brain: The role of uncertainty in neural coding and computation. *Trends in Neurosciences*, *27*(12), 712–719. https://doi.org/10.1016/j.tins.2004.10.007

Koblinger, Á., Fiser, J., & Lengyel, M. (2021). Representations of uncertainty: Where art thou? *Current Opinion in Behavioral Sciences*, *38*, 150–162. https://doi.org/10.1016/j.cobeha.2021.03.009

Körding, K. P., Beierholm, U., Ma, W. J., Quartz, S., Tenenbaum, J. B., & Shams, L. (2007). Causal inference in multisensory perception. *PLoS ONE*, *2*(9). https://doi.org/10.1371/journal.pone.0000943

Kristjánsson, Á. (2022). Priming of probabilistic attentional templates. *Psychonomic Bulletin & Review*. https://doi.org/10.3758/s13423-022-02125-w

Kristjánsson, Á., & Campana, G. (2010). Where perception meets memory: A review of repetition priming in visual search tasks. *Attention, Perception, & Psychophysics*, *72*(1), 5–18. https://doi.org/10.3758/APP.72.1.5

Kristjánsson, Á., & Driver, J. (2008). Priming in visual search: Separating the effects of target repetition, distractor repetition and role-reversal. *Vision research*, *48*(10), 1217–32. https://doi.org/10.1016/j.visres.2008.02.007

Lamy, D. F., Antebi, C., Aviani, N., & Carmel, T. (2008). Priming of pop-out provides reliable measures of target activation and distractor inhibition in selective attention. *Vision research*, *48*(1), 30–41. https://doi.org/10.1016/j.visres.2007.10.009

Landy, M. S., Banks, M. S., & Knill, D. C. (2011, September 14). Ideal-observer models of cue integration. In J. Trommershäuser, K. Kording, & M. S. Landy (Eds.), *Sensory cue integration* (pp. 5–29). Oxford University Press. https://doi.org/10.1093/acprof:oso/9780195387247.003.0001

Lange, R. D., Shivkumar, S., Chattoraj, A., & Haefner, R. M. (2020). Bayesian encoding and decoding as distinct perspectives on neural coding. *bioRxiv*, 1–16. https://doi.org/10.1101/2020.10.14.339770

Maljkovic, V., & Nakayama, K. (1994). Priming of pop-out: I. role of features. *Memory & cognition*, *22*(6), 657–72.

Manassi, M., Sayim, B., & Herzog, M. H. (2012). Grouping, pooling, and when bigger is better in visual crowding. *Journal of Vision*, *12*(10), 13. https://doi.org/10.1167/12.10.13

Oliva, A., & Torralba, A. (2001). Modeling the shape of the scene: A holistic representation of the spatial envelope. *International Journal of Computer Vision*, *42*(3), 145–175. https://doi.org/10.1023/A:1011139631724

Orhan, A. E., & Ma, W. J. (2015). Neural population coding of multiple stimuli. *Journal of Neuroscience*, *35*(9), 3825–3841. https://doi.org/10.1523/JNEUROSCI.4097-14.2015

Oriet, C., & Brand, J. (2013). Size averaging of irrelevant stimuli cannot be prevented. *Vision Research*, *79*, 8–16. https://doi.org/10.1016/j.visres.2012.12.004

Pascucci, D., Ceylan, G., & Kristjánsson, Á. (2022). Feature distribution learning by passive exposure. *Cognition*, *227*, 105211. https://doi.org/10.1016/j.cognition.2022.105211

Pewsey, A. (2004). The large-sample joint distribution of key circular statistics. *Metrika*, *60*(1). https://doi.org/10.1007/s001840300294

Portilla, J., & Simoncelli, E. P. (2000). Parametric texture model based on joint statistics of complex wavelet coefficients. *International Journal of Computer Vision*, *40*(1), 49–71. https://doi.org/10.1023/A:1026553619983

Pouget, A., Dayan, P., & Zemel, R. S. (2000). Information processing with population codes. *Nature Reviews Neuroscience*, *1*(2), 125–32. https://doi.org/10.1038/35039062

Qiu, C., Kersten, D., & Olman, C. A. (2013). Segmentation decreases the magnitude of the tilt illusion. *Journal of Vision*, *13*(13), 1–17. https://doi.org/10.1167/13.13.19

Rafiei, M., Chetverikov, A., Hansmann-Roth, S., & Kristjánsson, Á. (2021). You see what you look for: Targets and distractors in visual search can cause opposing serial dependencies. *Journal of Vision*, *21*(10), 3. https://doi.org/10.1167/jov.21.10.3

Rafiei, M., Hansmann-Roth, S., Whitney, D., Kristjánsson, Á., & Chetverikov, A. (2021). Optimizing perception: Attended and ignored stimuli create opposing perceptual biases. *Attention, Perception, & Psychophysics*, *83*(3), 1230–1239. https://doi.org/10.3758/s13414-020-02030-1

Rahnev, D. (2017). The case against full probability distributions in perceptual decision making. *bioRxiv*. https://doi.org/10.1101/108944

Rao, R. P., Olshausen, B. A., & Lewicki, M. S. (2002). *Probabilistic models of the brain: Perception and neural function*. MIT Press. https://doi.org/10.7551/mitpress/5583.001.0001

Rosenholtz, R. (2016). Capabilities and limitations of peripheral vision. *Annual Review of Vision Science*, *2*(1), 437–457. https://doi.org/10.1146/annurev-vision-082114-035733

Rosenholtz, R. (2020). Demystifying visual awareness: Peripheral encoding plus limited decision complexity resolve the paradox of rich visual experience and curious perceptual failures. *Attention, Perception, & Psychophysics*, *82*(3), 901–925. https://doi.org/10.3758/s13414-019-01968-1

Sahani, M., & Dayan, P. (2003). Doubly distributional population codes: Simultaneous representation of uncertainty and multiplicity. *Neural Computation*, *15*(10), 2255–2279. https://doi.org/10.1162/089976603322362356

Schwartz, O., Hsu, A., & Dayan, P. (2007). Space and time in visual context. *Nature Reviews Neuroscience*, *8*(7), 522–535. https://doi.org/10.1038/nrn2155

Seriès, P., & Seitz, A. R. (2013). Learning what to expect (in visual perception). *Frontiers in Human Neuroscience*, *7*(October), 668. https://doi.org/10.3389/fnhum.2013.00668

Shams, L., & Beierholm, U. R. (2010). Causal inference in perception. *Trends in Cognitive Sciences*, *14*(9), 425–432. https://doi.org/10.1016/j.tics.2010.07.001

Shurygina, O., Kristjánsson, Á., Tudge, L., & Chetverikov, A. (2019). Expectations and perceptual priming in a visual search task: Evidence from eye movements and behavior. *Journal of Experimental Psychology: Human Perception and Performance*, *45*(4), 489–499. https://doi.org/10.1037/xhp0000618

Stocker, A. A., & Simoncelli, E. P. (2005). Sensory adaptation within a bayesian framework for perception. *Advances in Neural Information Processing Systems*, 1289–1296.

Tanrıkulu, Ö. D., Chetverikov, A., Hansmann-Roth, S., & Kristjánsson, Á. (2021). What kind of empirical evidence is needed for probabilistic mental representations? an example from visual perception. *Cognition*, *217*, 104903. https://doi.org/10.1016/j.cognition.2021.104903

Tanrıkulu, Ö. D., Chetverikov, A., & Kristjánsson, Á. (2020). Encoding perceptual ensembles during visual search in peripheral vision. *Journal of Vision*, *20*(8), 20. https://doi.org/10.1167/jov.20.8.20

Tanrıkulu, Ö. D., Chetverikov, A., & Kristjánsson, Á. (2021). Testing temporal integration of feature probability distributions using role-reversal effects in visual search. *Vision Research*, *188*(July), 211–226. https://doi.org/10.1016/j.visres.2021.07.012

Treisman, A. (1996). The binding problem. *Current Opinion in Neurobiology*, *6*(2), 171–178. https://doi.org/10.1016/S0959-4388(96)80070-5

Treisman, A. (2006). How the deployment of attention determines what we see. *Visual Cognition*, *14*(4-8), 411–443. https://doi.org/10.1080/13506280500195250

Treisman, A., & Schmidt, H. (1982). Illusory conjunctions in the perception of objects. *Cognitive Psychology*, *14*(1), 107–141. https://doi.org/10.1016/0010-0285(82)90006-8

Utochkin, I. S., Khvostov, V. A., & Stakina, Y. M. (2018). Continuous to discrete: Ensemble-based segmentation in the perception of multiple feature conjunctions. *Cognition*, *179*(December 2017), 178–191. https://doi.org/10.1016/j.cognition.2018.06.016

Utochkin, I. S., & Vostrikov, K. O. (2017). The numerosity and mean size of multiple objects are perceived independently and in parallel. *PLoS ONE*, *12*(9), 1–20. https://doi.org/10.1371/journal.pone.0185452

Vértes, E., & Sahani, M. (2018). Flexible and accurate inference and learning for deep generative models. *Advances in Neural Information Processing Systems*, *2018-Decem*(NeurIPS), 4166–4175.

Wallis, T. S., Bethge, M., & Wichmann, F. A. (2016). Testing models of peripheral encoding using metamerism in an oddity paradigm. *Journal of Vision*, *16*(2), 1–30. https://doi.org/10.1167/16.2.4

Wei, X.-X., & Stocker, A. A. (2012). Efficient coding provides a direct link between prior and likelihood in perceptual bayesian inference. *Advances in Neural Information Processing Systems*, *25*(May), 1–9.

Wei, X.-X., & Stocker, A. A. (2015). A bayesian observer model constrained by efficient coding can explain 'anti-bayesian' percepts. *Nature Neuroscience*, *18*(10), 1509–1517. https://doi.org/10.1038/nn.4105

Whitney, D., & Yamanashi Leib, A. (2018). Ensemble perception. *Annual Review of Psychology*, *69*(1), 105–129. https://doi.org/10.1146/annurev-psych-010416-044232

Wolfe, J. M., Klempen, N. L., & Shulman, E. P. (1999). Which end is up? two representations of orientation in visual search. *Vision Research*, *39*(12), 2075–2086. https://doi.org/10.1016/S0042-6989(98)00260-0

Yeon, J., & Rahnev, D. (2020). The suboptimality of perceptual decision making with multiple alternatives. *Nature Communications*, *11*(1), 1–12. https://doi.org/10.1038/s41467-020-17661-z

Zemel, R. S., Dayan, P., & Pouget, A. (1998). Probabilistic interpretation of population codes. *Neural Computation*, *10*(2), 403–430. https://doi.org/10.1162/089976698300017818

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

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

## | Task structure

Participants have to locate a target among a set of distractors and indicate if it is in the top or in the lower part of the stimuli matrix. The experimenter sets the task parameters, namely, the target distribution, $p(s_i | L_T = i)$, and the distractor distribution, $p(s_i | L_T \neq i)$, for each location $i = 1 \ldots N$ in the stimuli matrix (with top half having indices from 1 to $N/2$ and the 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 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, that the target location is or is not $i$, respectively.

## | Ideal observer model

At each moment in time $t = 1 \ldots K$ (with $K$ as the decision moment) and at each location $i$, the observer obtains sensory observations $x_{i,t}$ corrupted by the presence of sensory noise:

$$p(x_{i,t} \mid s_i) = f_{VM}(x_i; s_i, \kappa_s) \qquad (S.1)$$

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

$$p(\mathbf{X} \mid \mathbf{s}) = \prod_{i=1}^{N} p(\mathbf{x_i} \mid s_i) = \prod_{i=1}^{N} \prod_{t=1}^{K} p(x_{i,t} \mid s_i) \qquad (S.2)$$

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

$$d = \frac{p(C = 1 \mid \mathbf{X})}{p(C = 2 \mid \mathbf{X})} \qquad (S.3)$$

where $C = 1$ and $C = 2$ correspond to the two hypotheses about the target location. After applying the log transfor-

mation, the decision variable can be expressed as a difference in the amount of evidence for the two hypotheses:

$$\log d = \log p\left(C = 1 \mid \mathbf{X}\right) - \log p\left(C = 2 \mid \mathbf{X}\right) \tag{S.4}$$

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

$$K = \frac{\log B}{E\left[\log p\left(C = 1 \mid \mathbf{X}\right)\right] - E\left[\log p\left(C = 2 \mid \mathbf{X}\right)\right]} \tag{S.5}$$

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

$$p\left(C = 1 \mid \mathbf{X}\right) = \frac{p\left(\mathbf{X} \mid C = 1\right)p(C = 1)}{p\left(\mathbf{X}\right)} \tag{S.6}$$

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

$$p\left(\mathbf{X} \mid C = 1\right) = \int p\left(\mathbf{X} \mid \mathbf{s}\right)p\left(\mathbf{s} \mid C = 1\right)d\mathbf{s} \tag{S.7}$$

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

$$p\left(\mathbf{X} \mid C = 1\right) = \sum_{i=1}^{\frac{N}{2}} \int p\left(\mathbf{X} \mid \mathbf{s}\right)p^*\left(\mathbf{s} \mid L_T = i, \boldsymbol{\theta}\right)d\mathbf{s} \tag{S.8}$$

where similarly to the main text, we use an asterisk to denote probability distributions as approximated by the observer through a set of parameters related to previous and current trials $\boldsymbol{\theta} = \{\boldsymbol{\theta}_{prev}, \boldsymbol{\theta}_{curr}\}$. That is, we assume that the observer is unaware of the true distributions $p(s_i | L_T = i)$ and $p(s_i | L_T \neq i)$ and approximates them instead using the information available. Note that the sum is done separately for each half of the stimuli matrix, hence $\frac{N}{2}$ is used in Eq. S.8.

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

$$p^*\left(\mathbf{s} \mid L_T = i, \boldsymbol{\theta}\right) = p^*\left(s_i \mid L_T = i, \boldsymbol{\theta}\right)\prod_{j \neq i}^{N} p^*\left(s_j \mid L_T \neq j, \boldsymbol{\theta}\right) \tag{S.9}$$

Using Eq. S.9, it can be further shown that:

$$\int p\left(\mathbf{X} \mid \mathbf{s}\right)p^*\left(\mathbf{s} \mid L_T = i, \boldsymbol{\theta}\right)d\mathbf{s} = \left[\prod_{j}^{N} \int p\left(\mathbf{x}_j \mid s_j\right)p^*\left(s_j \mid L_T \neq j, \boldsymbol{\theta}\right)ds_j\right]\frac{\int p\left(\mathbf{x}_i \mid s_i\right)p^*\left(s_i \mid L_T = i, \boldsymbol{\theta}\right)ds_i}{\int p\left(\mathbf{x}_i \mid s_i\right)p^*\left(s_i \mid L_T \neq i, \boldsymbol{\theta}\right)ds_i} \tag{S.10}$$

Note that the product in the square brackets is the same for all locations, and the remaining part of the equation is a

ratio of the probability that the measurements at a given location are from the target against the probability that they are from the distractor, similarly to the model described in the main text.

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

$$p^* \left( s_i \mid L_T = i, \boldsymbol{\theta} \right) = p^* \left( s_i \mid L_T = i, \boldsymbol{\theta}_{prev} \right) p^* \left( s_i \mid L_T = i, \boldsymbol{\theta}_{curr} \right) \tag{S.11}$$

For each location and each location-specific hypothesis $L_T = i$ and $L_T \neq i$, the current trial parameters need to be computed separately because of the nature of the odd-one-out task. A target is defined as the item most different from the distractors. For simplicity, we assumed that observers use the following circular normal approximation for the distractors at the current trial based on the sensory observations:

$$p^* \left( s_i \mid L_T \neq i, \boldsymbol{\theta}_{curr} \right) = f_{VM} \left( s_i; \hat{\mu}_{j \neq i}, \hat{\kappa}_{j \neq i} \right) \tag{S.12}$$

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

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

$$p^* \left( s_i \mid L_T = i, \boldsymbol{\theta}_{curr} \right) = f_{VM} \left( s_i; \hat{\mu}_{j \neq i} + 90°, \kappa_T \right) \tag{S.13}$$

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

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

$$p \left( \mathbf{X} \mid C = 1 \right) = \left[ \prod_{j=1}^{N} \int \left( \mathbf{x}_j \mid s_j \right) p^* \left( s_j \mid L_T \neq j, \boldsymbol{\theta} \right) ds_j \right] \sum_{i=1}^{\frac{N}{2}} \frac{\int p \left( \mathbf{x}_i \mid s_i \right) p^* \left( s_i \mid L_T = i, \boldsymbol{\theta} \right) ds_i}{\int p \left( \mathbf{x}_i \mid s_i \right) p^* \left( s_i \mid L_T \neq i, \boldsymbol{\theta} \right) ds_i} \tag{S.14}$$

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

$$\log d = \log \left( \sum_{i=1}^{\frac{N}{2}} \frac{\int p \left( \mathbf{x}_i \mid s_i \right) p^* \left( s_i \mid L_T = i, \boldsymbol{\theta} \right) ds_i}{\int p \left( \mathbf{x}_i \mid s_i \right) p^* \left( s_i \mid L_T \neq i, \boldsymbol{\theta} \right) ds_i} \right) - \log \left( \sum_{i=\frac{N}{2}+1}^{N} \frac{\int p \left( \mathbf{x}_i \mid s_i \right) p^* \left( s_i \mid L_T = i, \boldsymbol{\theta} \right) ds_i}{\int p \left( \mathbf{x}_i \mid s_i \right) p^* \left( s_i \mid L_T \neq i, \boldsymbol{\theta} \right) ds_i} \right) \tag{S.15}$$

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

## | Simulations

To estimate the behavior of the observer using this model, we simulated the decision-making process and estimated the mean response times while varying the properties of the distractor representation $p^*\left(s_i \mid L_T \neq i, \theta_{prev}\right)$. The task parameters were based on the actual experiment design. We used 36 stimuli for each trial with one stimulus being the test target ($s_{L_T}$) and the rest being the distractors. The distractors on each simulated trial were distributed as $p\left(s_i \mid L_T \neq i\right) = f_{VM}\left(s_i; \mu_D, \kappa_D\right)$ where $\mu_D \sim U\left(s_{L_T} + 60°; s_{L_T} + 120°\right)$ (that is, the mean of distractors is set to 60° to 120° away from the test stimulus) and $\kappa_D = 8.7$ (approximately equivalent to the standard deviation of 10° in orientation space). The sensory observations were assumed to be noisy ($\kappa_s = 2$, approximately equivalent to the standard deviation of 24° in orientation space; note that this is the noise level for samples collected at each moment in time). The observers' target representation was assumed to be linked with to the distractor representation as $p^*\left(s_i \mid L_T = i, \theta_{prev}\right) = f_{VM}\left(s_i; \mu_{D_{prev}}, \kappa_T\right)$ with $\kappa_T = 3.35$ (based on a normal approximation to a uniform target distribution with 60° range used in the experiments). The same $\kappa_T$ was used for target-related computations based on the current trial data (Eq. S.13). The decision threshold was set to $\log B = 4.60$ assuming a 1% probability of error if the observer assumptions are correct. For each test target from 1° to 180° in half-degree steps, we simulated 56 trials for each combination of distractor representation parameters.

We ran simulations for the wrapped skewed normal distribution with the mean varied from -60° to 60° in 20° steps, while the standard deviation varied from 20° to 60° in 10° steps, and skew varied from -10 to 10 in steps of 2. The results of the simulations (Figure S3) confirmed the findings obtained with a simplified model: the means are recovered precisely while for standard deviation and skewness the monotonic relationship holds.
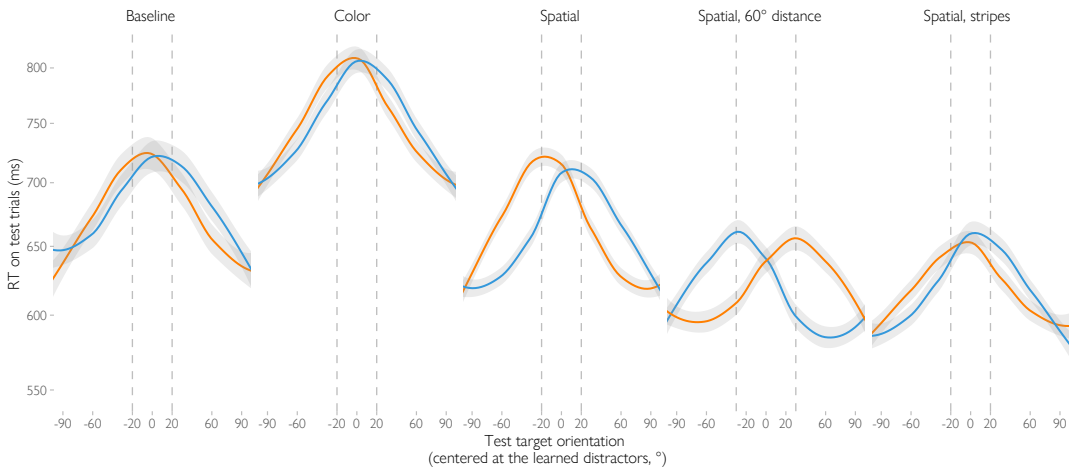
**FIGURE S1** Raw response times for test trials in different conditions. The curves shows the average RT with 95% CI shown with shading. Different colors indicate RT for test trials when a test target matches (in location or color) the clockwise- and counterclockwise-shifted parts of the learned distractor distribution. Dashed lines show the mean orientations for the corresponding distribution parts. The average was estimated with locally-weighted regression (LOESS, an extension of the moving window approach) that accounted for circularity of the orientation space by padding the data at each end of the [-90,90] range with 1/6$^{th}$ of the data from the other end.
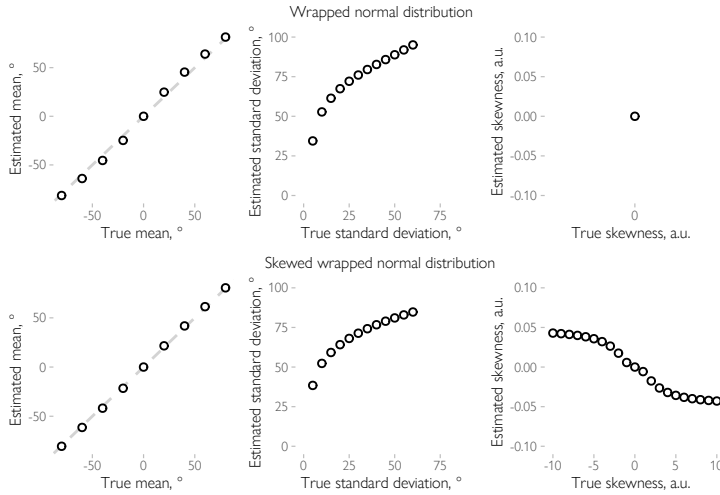
**FIGURE S2** Simulated parameters under the simplified Bayesian observer model. We simulated the response times under the assumptions of the simplified Bayesian observer model described in the main text and applied the same approach as used for the real data to see if the assumed monotonic relationship between the true parameters and the recovered parameters holds. Firstly, we used a simple wrapped normal (top) with means varying from -80° to 80° in 20° steps and standard deviation from 5° to 60° in 5° steps. For each parameter combination the RT were computed using Eq. 16. 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, 2004) with means again varying from -80° to 80° in 20° steps, scale parameter varying from 5° to 60° in 5° 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.
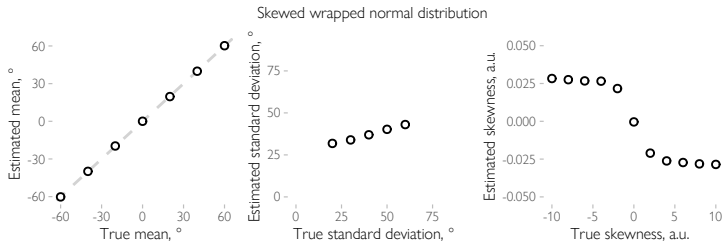
Skewed wrapped normal distribution

**FIGURE S3** 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 S2). 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.