

A rational theory of the limitations of working memory and attention

Ronald van den Berg¹ and Wei Ji Ma²

¹Department of Psychology, University of Uppsala, Uppsala, Sweden.

²Center for Neural Science and Department of Psychology, New York University, New York, USA

Corresponding author

Ronald van den Berg

Department of Psychology

Von Kraemers Allé 1A, 75237

Uppsala, Sweden

Tel: +46184717277

E-mail: ronald.vandenberg@psyk.uu.se

The precision with which items are encoded in working memory and attention decreases with the number of encoded items¹⁻⁴. Current theories typically account for this “set size effect” by postulating a hard constraint on the allocated amount of encoding resource, commonly formalized as samples, spikes, slots, or bits. While these theories have produced models that are quantitatively successful, they offer no principled explanation for the very existence of set size effects: given their detrimental consequences for behavioral performance, why have these effects not been weeded out by evolutionary pressure, for example by scaling the amount of allocated encoding resource with set size? Here, we propose a theory that is based on an ecological notion of rationality: set size effects establish an optimal trade-off between behavioral performance and the neural costs associated with stimulus encoding. We derive models from this theory for four visual working memory and attention tasks and find that it accounts well for data in eleven different experiments. Our results suggest that set size effects have a rational basis and that ecological costs should be considered in models of human behavior.

Set size effects in working memory and attention have been modeled extensively, but no model offers a principled explanation for the existence of such effects. Models typically postulate that stimuli are encoded using a fixed total amount of resources, formalized as “samples”³⁻⁵, slots⁶, information bit rate⁷, Fisher information⁸, or neural firing⁹. Apart from quantitative deviations from the data^{10,11}, these models do not explain why the brain would not use more resources when set size increases. Another class of models does allow for total resources to be set size dependent, for example by assuming a power law relationship between precision and set size^{2,11-17}; these models tend to provide excellent fits, but have been criticized to lack a principled explanation of set size effects^{18,19}.

Here, we derive models for set size effects starting from the ecological principle that neural firing is energetically costly²⁰⁻²². This cost may have pressured the brain to balance behavioral benefits of high precision against the neural loss that it induces^{8,21,23,24}. What level of encoding precision establishes a good balance might depend on multiple factors, such as set size, task, and motivation. Indeed, performance on perceptual decision-making tasks can be improved by increasing monetary reward²⁵⁻²⁷, which suggests that the total amount of resource spent on encoding has some flexibility that is driven by ecological factors. Based on these considerations, we hypothesize that set size effects on encoding precision reflect an ecologically rational strategy

that balances behavioral performance against neural costs. Next, we derive formal models from this hypothesis and test it on four visual working memory and attention tasks.

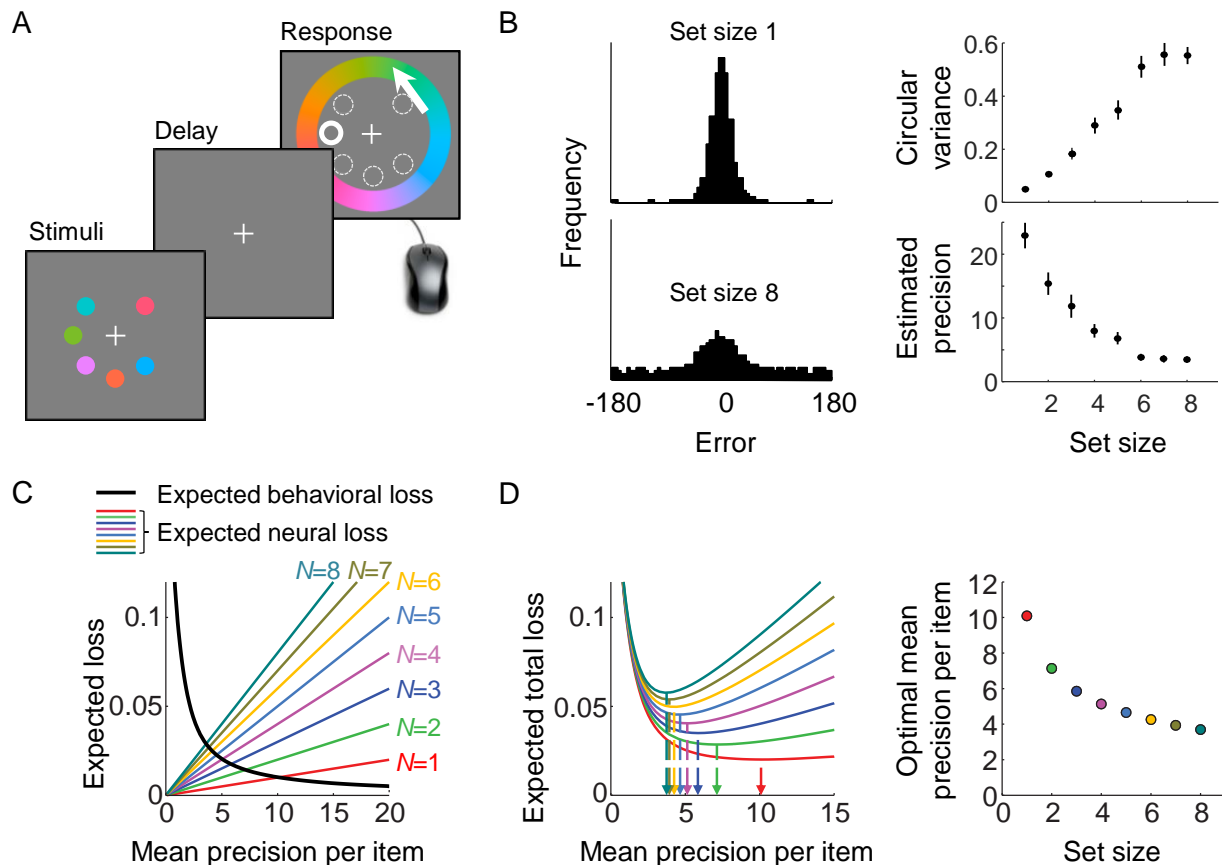


Figure 1. An ecologically rational model of set size effects in delayed estimation. (A) Example of a delayed-estimation experiment. The subject is briefly presented with a set of stimuli and, after a short delay, reports the value of a randomly chosen target item. (B) Estimation error distributions widen with set size, suggesting a decrease in encoding precision (data from Experiment DE5 in Table 1; estimated precision computed in the same way as in Fig. 2C). (C) Stimulus encoding is assumed to be associated with two kinds of loss: a behavioral loss that decreases with encoding precision and a neural loss that is proportional to both set size and precision. In the delayed-estimation task, the expected behavioral error loss is independent of set size. (D) Total expected loss has a unique minimum that depends on the number of remembered items. The mean precision per item that minimizes expected total loss is referred to as the optimal mean precision (arrows) and decreases with set size. The parameter values used to produce panels C and D were $\lambda=0.01$, $\beta=2$, and $\tau \downarrow 0$.

Table 1. Overview of experimental datasets used. Task responses were continuous in the delayed-estimation experiments and categorical in the other tasks. DE5 and DE6 differed in the way color was reported (DE5: color wheel; DE6: scroll).

Experiment	Reference	Task	Feature	Set sizes	#subj
DE1	¹⁷	Delayed estimation	Color	1, 2, 4, 8	15
DE2	⁶	Delayed estimation	Color	1, 2, 3, 6	8
DE3	¹²	Delayed estimation	Color	1, 2, 4, 6	12
DE4	¹⁴	Delayed estimation	Orientation	1-8	6
DE5	¹⁴	Delayed estimation	Color	1-8	13
DE6	¹⁴	Delayed estimation	Color	1-8	13
CD1	¹⁰	Change detection	Color	1, 2, 4, 8	7
CD2	¹⁰	Change detection	Orientation	2, 4, 6, 8	10
CL1	¹⁴	Change localization	Color	2, 4, 6, 8	7
CL2	¹⁴	Change localization	Orientation	2, 4, 6, 8	11
VS	²⁸	Visual search	Orientation	1, 2, 4, 8	6

The first task that we consider is delayed estimation ¹⁷, in which the subject briefly holds a set of items in memory and reports their estimate of a randomly chosen target item (Fig. 1A; Table 1). Estimation error ε is the circular difference between the subject's estimate and the true stimulus value s . Set size effects in this task are visible as a widening of the estimation error distribution (Fig. 1B). As in previous work ^{2,10,11,14,28,29}, we assume that a memory x follows a Von Mises distribution with mean s and concentration parameter κ , and define precision J as Fisher information ³⁰. J is one-to-one related to κ (see Supplementary Information), and in the absence of response noise, $\varepsilon=x-s$. Moreover, we assume variability in J across items and trials ^{2,11,14,29,31}, which we model using a gamma distribution with a mean \bar{J} and a scale parameter τ (see Supplementary Information).

The key idea of our theory is that stimuli are encoded with a level of (mean) precision, \bar{J} , that minimizes the combination of two kinds of loss. The first one is the *behavioral loss* induced by making an error ε , as described by a mapping $L_{\text{behavioral}}(\varepsilon)$, which may depend on both internal incentives (e.g., intrinsic motivation) and external ones (e.g., the reward scheme imposed by the experimenter). For the moment, we choose a power-law function, $L_{\text{behavioral}}(\varepsilon)=|\varepsilon|^\beta$, with

$\beta > 0$ as a free parameter. The *expected* behavioral loss, denoted $\bar{L}_{\text{behavioral}}$, is obtained by averaging loss across all possible errors, weighted by the probability that each error occurs,

$$\bar{L}_{\text{behavioral}}(\bar{J}, N) = \int L_{\text{behavioral}}(\varepsilon) p(\varepsilon | \bar{J}, N) d\varepsilon, \quad (1)$$

where $p(\varepsilon | \bar{J}, N)$ is the estimation error distribution for given mean precision and set size. In single-probe delayed-estimation tasks, the expected behavioral loss is independent of set size and subject to the law of diminishing returns (Fig. 1C, black curve).

The second kind of loss is a *neural loss* induced by the neural spiking activity that represents a stimulus. For many choices of spike variability, including the common choice of Poisson-like variability³², the precision (Fisher information) of a stimulus encoded in a neural population is proportional to the neural gain^{33,34}. We assume for the moment that the neural loss is proportional to neural gain and, therefore, to precision. Further assuming that stimuli are encoded independently of each other, neural loss is also proportional to the number of encoded items, N . We thus obtain

$$\bar{L}_{\text{neural}}(\bar{J}, N) = \alpha \bar{J} N, \quad (2)$$

where α is a free parameter that represents the amount of neural loss incurred by a unit increase in mean precision (Fig. 1C, colored lines).

We combine the two types of expected loss into a total expected loss function (Fig. 1D) by taking a weighted sum,

$$\begin{aligned} \bar{L}_{\text{total}}(\bar{J}, N) &= \bar{L}_{\text{behavioral}}(\bar{J}, N) + \lambda \bar{L}_{\text{neural}}(\bar{J}, N) \\ &= \bar{L}_{\text{behavioral}}(\bar{J}, N) + \lambda \alpha \bar{J} N, \end{aligned} \quad (3)$$

where the weight λ represents the importance of keeping neural loss low relative to the importance of good performance. Since λ and α have interchangeable effects on the model predictions, they can be fitted as a single free parameter $\tilde{\lambda} \equiv \lambda \alpha$. We refer to the level of mean precision that minimizes the total expected loss as *optimal mean precision* (Fig. 1D),

$$\bar{J}_{\text{optimal}}(\bar{J}, N) = \underset{J}{\operatorname{argmin}} \bar{L}_{\text{total}}(\bar{J}, N). \quad (4)$$

It can be proven mathematically (see Supplementary Information) that this model predicts a decrease in mean encoding precision with set size if the following four conditions are satisfied: (i) expected behavioral loss is a strictly decreasing function of encoding precision, i.e., an increase in precision results in an increase in performance; (ii) expected behavioral loss is subject to a law of diminishing returns³⁵: the higher the initial precision, the smaller the behavioral benefit obtained from an increase in precision; (iii) expected neural loss is an increasing function of encoding precision; (iv) expected neural loss is either constant or increases with encoding precision.

We used maximum-likelihood estimation to fit the three model parameters ($\tilde{\lambda}$, τ , and β) to 67 individual-subject data sets from a delayed-estimation benchmark set* (Table 1). The model accounts well for the raw error distributions and the two statistics that summarize these distributions (Fig. 2A). Model comparison based on the Akaike Information Criterion (AIC)³⁶ indicates that the goodness of fit is comparable to that of a descriptive model variant in which the relation between encoding precision and set size is assumed to follow a power law ($\Delta\text{AIC}=5.27\pm0.70$ in favor of the rational model). Hence, the rational model provides a principled explanation of set size effects in delayed-estimation tasks without sacrificing quality of fit.

We next try to falsify our theory by testing whether a mapping between set size and encoding precision exists that fits the data better than the relation imposed by the loss-minimization strategy of the rational model. To this end, we fit an unconstrained variant of the model in which memory precision is fitted as a free parameter at each set size. We find only a minimal difference in goodness of fit ($\Delta\text{AIC}=3.49\pm0.93$ in favor of the unconstrained model), suggesting that the fits of the rational model are close to the best possible fits. This finding is corroborated by examination of the fitted parameter values: the estimated precision values in the unconstrained model closely match the precision values in the rational model (Fig. 2B). Hence, it seems that no relation exists that fits these data substantially better than the set of relations that the rational model is constrained to.

* The original benchmark set¹¹ contains 10 data sets with a total of 164 individuals. Two of these data sets were published in papers that later got retracted and another one contained data for only two set sizes, which is not very informative for our present purposes. While our model accounts well for these data sets (Fig. S1 in Supplementary Information), we decided to exclude them from the main analyses.

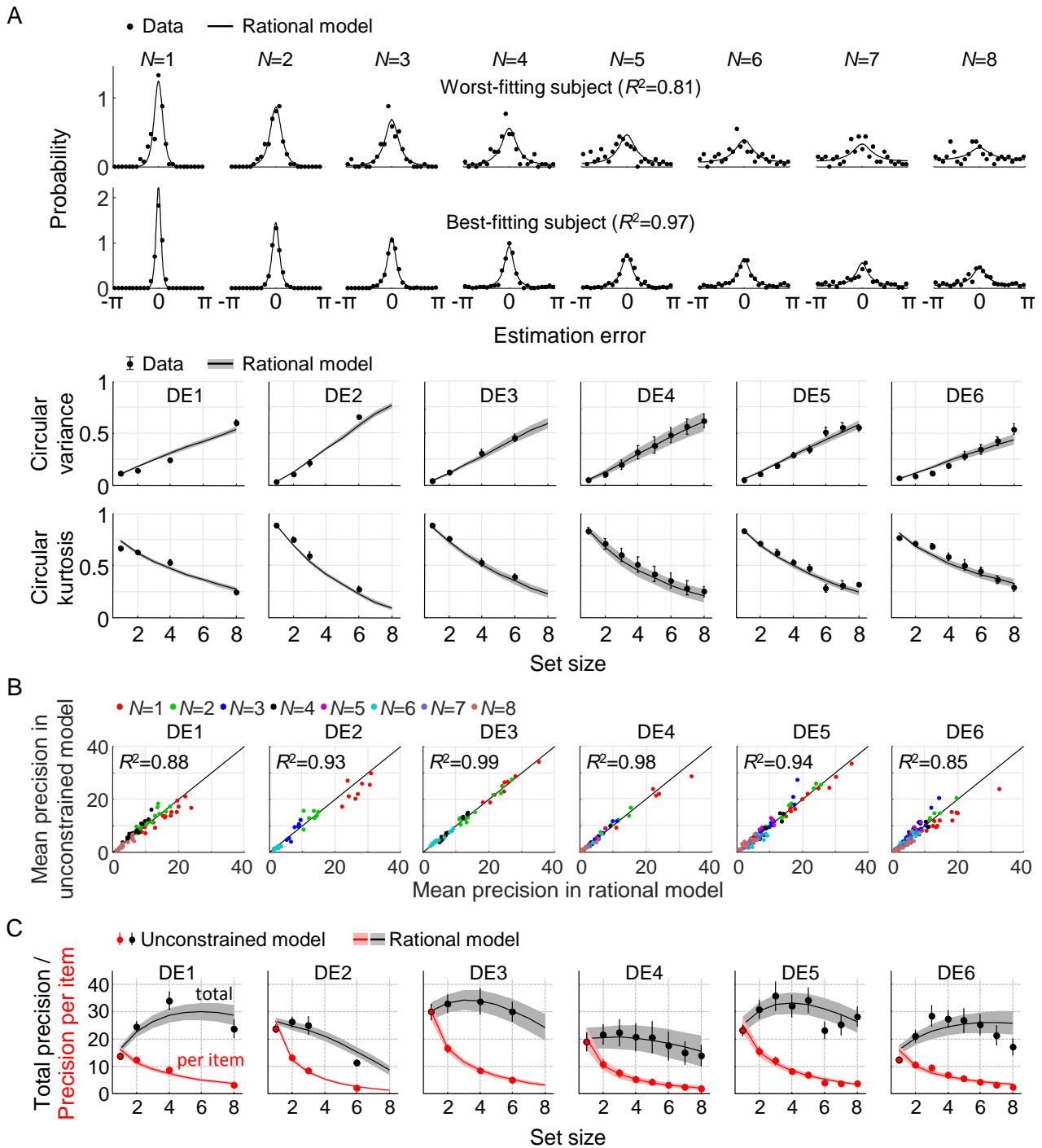


Figure 2. Model fits to data from six delayed-estimation experiments. (A) Maximum-likelihood fits to raw data of the worst-fitting and best-fitting subjects (based on R^2). (B) Subject-averaged fits to the two statistics that summarize the estimation error distributions (circular variance kurtosis) as a function of set size, split by experiment. Here and in subsequent figures, error bars and shaded areas represent 1 s.e.m. of the mean across subjects. (C) Best-fitting precision values in the rational model scattered against the best-fitting precision values in the unconstrained model. Each dot represents the estimate for a single subject. (D) Estimated mean encoding precision per item (red) and total encoding precision (black) plotted against set size.

Further analysis of the estimated precision values in the unconstrained model suggests that total memory precision (defined as $\bar{J}_{\text{total}} = \bar{J}N$) varies non-monotonically with set size (Fig. 2C, black circles). Whereas models that assume a power-law relation between precision per item and set size constrain \bar{J}_{total} to be a monotonic function of N , non-monotonic relations are predicted by the rational model (Fig. 2C, gray curves).

To evaluate the necessity of a free parameter in the behavioral loss function, $L_{\text{behavioral}}(\varepsilon)$, we also test the following three parameter-free choices: $|\varepsilon|$, ε^2 , and $-\cos(\varepsilon)$. Model comparison favors the original model with AIC differences of 14.0 ± 2.8 , 24.4 ± 4.1 , and 19.5 ± 3.5 , respectively. While there may be other parameter-free functions that give better fits, we expect that a free parameter is unavoidable here, as it is likely that the error-to-loss mapping differs across experiments (due to differences in external incentives) and possibly also across subjects within an experiment (due to differences in internal incentives). We also test a two-parameter function that was proposed recently (Eq. (5) in ³⁷). The main difference with our original choice is that this alternative function allows for saturation effects in the error-to-loss mapping. However, this extra flexibility does not increase the goodness of fit, as the original model outperforms this variant with an AIC difference of 5.3 ± 1.8 .

We next examine the generality of our theory, by testing whether it can also explain set size effects in two change detection tasks (Table 1). In these experiments, the subject is on each trial sequentially presented with two sets of stimuli and reports whether there was a change at any of the stimulus locations (Fig. 3A). A change was present on half of the trials, at a random location and with a random change magnitude. The behavioral error, ε , takes only two values in this task: “correct” and “incorrect”. Therefore, $p(\varepsilon | \bar{J}, N)$ specifies the probabilities of correct and incorrect responses for a given level of precision and set size, which depend on the observer’s decision rule. Following previous work ^{10,29}, we assume that subjects use the Bayes-optimal decision rule (see Supplementary Information) and that there is random variability in encoding precision. This decision rule introduces one free parameter, p_{change} , which specifies the subject’s degree of prior belief that a change will occur. Due to the binary nature of ε in this task, the free parameter of the behavioral loss function drops out of the model, because its effect is equivalent to changing parameter $\tilde{\lambda}$ (see Supplementary Information). The rational model thus

has three free parameters ($\tilde{\lambda}$, τ , and p_{change}). We find that the maximum-likelihood fits account well for the data in both experiments (Fig. 3B).

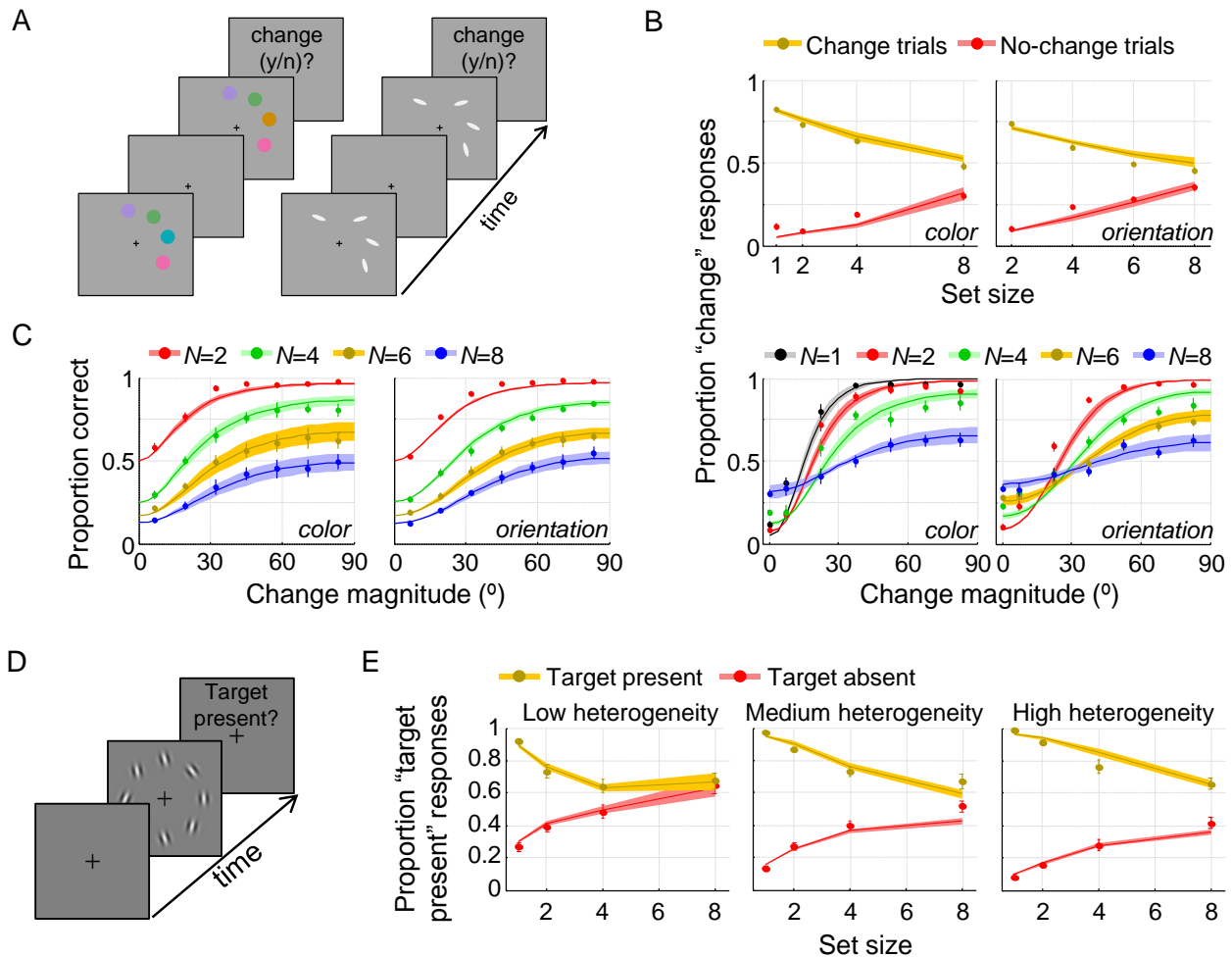


Figure 3. Model fits to three categorical decision-making tasks. (A) Experimental paradigm in the change-detection experiments. The paradigm for change localization was the same, except that a change was present on each trial and subjects reported the location of change. (B) Model fits to change-detection data. Top: hit and false alarm rates; bottom: psychometric curves. (C) Model fits to change-localization data. (D) Experimental paradigm in the visual-search experiment. (E) Model fits to visual-search data.

So far, we have considered tasks with continuous and binary judgments. We next consider two change localization experiments (Table 1) in which judgments are non-binary but categorical. The task is identical to change detection, except that a change is present on every trial and the observer reports the location at which the change occurred (out of 2, 4, 6, or 8

locations). We again assume variable precision and an optimal decision rule (see Supplementary Information). Although the rational model has only two free parameters ($\tilde{\lambda}$ and τ), it accounts well for both datasets (Fig. 3C).

The final task to which we apply our theory is a visual search experiment²⁸ (Table 1). Unlike the previous three tasks, this is not a working memory task, as there was no delay period between stimulus offset and response. Set size effects in this experiment are thus likely to stem from limitations in attention rather than memory, but our theory applies without any additional assumptions. Subjects judged whether a vertical target was present among one of N briefly presented oriented ellipses (Fig. 3D). The distractors were drawn from a Von Mises distribution centered at vertical. The width of the distractor distribution determined the level of heterogeneity in the search display. Each subject was tested under three different levels of heterogeneity. We again assume variable precision and an optimal decision rule (see Supplementary Material). This decision rule has one free parameter, p_{present} , specifying the subject's prior degree of belief that a target will be present. We fit the three free parameters ($\tilde{\lambda}$, τ , and p_{present}) to the data from all three heterogeneity conditions at once. The rational model accounts well for the hit and false alarm rates as a function of set size and their dependence on heterogeneity condition (Fig. 3E).

A key strength of our theory is that it uses a single principle of rationality and relatively few parameters to produce well-fitting models across a range of quite different tasks. Nevertheless, consideration of additional mechanisms could further improve the fits and lead to more complete models of human behavior. For example, previous studies have incorporated response noise^{11,14}, non-target responses¹², and a (variable) limit on the number of remembered items^{7,11,38} to improve fits. We did not consider such components here, as they come with additional parameters, some are task-specific (such as non-target responses), and they have so far not been motivated in a principled manner. Regarding the latter point, it might be possible to treat some of these mechanisms using an ecologically rational approach as well. For example, the level of response noise might be set by optimizing a trade-off between performance and motor control effort³⁹.

Our work speaks to the relation between descriptive and rational theories in psychology and neuroscience. The main motivation for rational theories is to reach a deeper level of understanding by analyzing a system in the context of the ecological needs and constraints that it evolved under. Besides the large literature on ideal-observer decision rules^{40–43}, rational

approaches have been used to explain properties of receptive fields^{44–46}, tuning curves^{47–49}, neural wiring^{50,51}, and neural network modularity⁵². A transition from descriptive to rational explanations might be an essential step in the maturation of theories of biological systems, and in psychology there certainly seems to be more room for this kind of explanation.

Although several previous models in the field of working memory and attention contain rational aspects, none of them accounts for set size effects in a principled way. Sims and colleagues have examined how errors in visual working memory can be minimized by optimally taking into account statistics of the stimulus distribution, but assume a fixed total amount of available encoding resource^{7,53}. Moreover, in our own previous work on visual search^{2,28}, change detection^{10,29}, and change localization¹⁴, we used optimal-observer models for the decision stage, but assumed an ad hoc power law for the encoding stage. An alternative explanation of set size effects has been that the brain is unable to keep neural representations of multiple items segregated from one another^{19,54–56}: as the number of encoded items increases, so does the level of interference in their representations, resulting in lower task performance. However, these models offer no rational justification for the existence of interference and some require additional mechanisms to account for set size effects; for example, the model by Oberauer and colleagues requires three additional components – including a set-size dependent level of background noise – to fully account for set size effects¹⁹. That being said, our theory does not rule out the possibility of interference, and it could be added onto any of the models we presented.

Our approach shares both similarities and differences with the concept of bounded rationality⁵⁷, which states that human behavior is guided by mechanisms that provide “good enough” solutions rather than optimal ones. The main similarity is that both approaches acknowledge that human behavior is constrained by various cognitive limitations. However, an important difference is that bounded rationality postulates these limitations as a given fact, while our approach explains them as rational outcomes of ecological optimization processes. The suggestion that cognitive limitations are themselves subject to optimization may also have implications for theories outside the field of psychology. One example concerns recent models of value-based decision-making that incorporate constraints imposed by working memory and attention limitations (e.g.,⁵⁸). Another example is the theory of “rational inattention” in behavioral economics, which examines optimal decision-making under the assumption that

decision makers have a fixed limit on the total amount of attention that they can allocate to process economic data⁵⁹. It might be interesting to extend that theory by treating the amount of allocable attention as the outcome of an optimization process rather than a constant.

While our results show that set size effects can in principle be explained as the result of an optimization strategy, they do not necessarily imply that encoding precision is fully optimized on every trial at any given task. First, encoding precision in the brain most likely has an upper limit, due to irreducible sources of noise such as Johnson noise and Poisson shot noise⁶⁰, as well as suboptimalities early in sensory processing⁶¹. This prohibits subjects to reach the near-perfect performance levels that our model may predict when the behavioral loss associated to errors is huge. Second, precision might have a lower limit: task-irrelevant stimuli are sometimes automatically encoded⁶², perhaps because in natural environments few stimuli are ever completely irrelevant. This would prevent subjects from sometimes encoding nothing at all, which our theory predicts to happen at very large set sizes. Third, all models that we tested incorporated variability in encoding precision. Part of this variability is possibly due to stochastic factors such as neural noise, but part of it may also be systematic in nature (e.g., particular colors and orientations may be encoded with higher precision than others^{63,64}). Whereas the systematic component could have a rational basis (e.g., higher precision for colors and orientations that occur more frequently in natural scenes^{49,65}), this is unlikely to be true for the random component. Indeed, when we jointly optimize \bar{J} and τ , we find estimates of τ that are consistently 0, which corresponds to absence of variability.

Future work could further examine optimality of encoding precision in working memory and attention by studying effects of changing experimental parameters that affect the loss functions. In delayed estimation, an obvious choice for this would be the delay period. One possibility is that working memories are maintained in persistent activity^{66,67}, in which case a longer delay would induce a higher cost and decrease optimal encoding precision. However, the relationship between neural cost and delay is not known, which makes it difficult to make precise quantitative predictions. In addition, it has been argued that working memories may be stored through rapid changes in synaptic strength, without the need of enhanced spiking activity⁶⁸. In that case, a longer delay might not even induce a higher neural cost. Another experimental parameter that could be varied is the error-to-loss mapping. A previous study that performed this manipulation found an effect in one experiment, but did not vary set size⁶⁹. None of the

experiments modeled here contained this manipulation (DE4-DE6 imposed an explicit loss function but did not vary it; the other experiments had no explicit scoring system). Future studies could measure effects of changes in explicitly imposed scoring systems and test how well a rational model accounts for such effects. However, it is important to keep in mind that subjects may not be able to fully internalize experimental loss functions in the timespan of a single experiment. This would mean that there generally is a mismatch between the loss functions used by the subject and the ones that would produce an optimal trade-off between performance and neural loss, which calls for further caution when testing of model predictions.

Finally, our results raise the question what neural mechanisms could implement the kind of near-optimal resource allocation strategy that is the core of our theory. Some form of divisive normalization^{9,70} would be a likely candidate, as it has the effect of lowering the gain when set size is larger. Moreover, divisive normalization is already a key operation in neural models of attention⁷¹ and visual working memory^{9,54}.

METHODS

Data and code availability

All data analyzed in this paper and model fitting code are available at [url to be inserted].

Model fitting

Delayed estimation. We used Matlab's `fminsearch` function to find the parameter vector

$\theta = \{\tilde{\lambda}, \beta, \tau\}$ that maximizes the log likelihood function, $\sum_{i=1}^n \log p(\varepsilon_i | N_i, \theta)$, where n is the

number of trials in the subject's data set, ε_i the estimation error on the i^{th} trial, and N_i the set size on that trial. To reduce the risk of converging into a local maximum, initial parameter estimates were chosen based on a coarse grid search over a large range of parameter values. The predicted estimation error distribution for a given parameter vector θ was computed as follows. First, \bar{J}_{optimal} was computed by applying Matlab's `fminsearch` function to Eq. (5). In this process, the integrals over ε and J were approximated numerically by discretizing the distributions of these variables into 100 and 20 equal-probability bins, respectively. Next, the gamma distribution over precision with mean \bar{J}_{optimal} and scale parameter τ was discretized into 20 equal-probability bins.

Thereafter, the predicted estimation error distribution was computed under the central value of each bin. Finally, these 20 predicted distributions were averaged. We verified that our results are robust under changes in the number of bins used in the numerical approximations.

Change detection. Model fitting in the change detection task consisted of finding parameter vector $\theta = \{\tilde{\lambda}, \tau, p_{\text{change}}\}$ that maximizes $\sum_{i=1}^n \log p(R_i | \Delta_i, N_i, \theta)$, where n is the number of trials in the subject's data set, R_i is the response ("change" or "no change"), Δ_i the magnitude of change, and N_i the set size on the i^{th} trial. For computational convenience, Δ was discretized into 30 equally spaced bins. To find the maximum-likelihood parameters, we first created a table with predicted probabilities of "change" responses for a large range of $(\bar{J}, \tau, p_{\text{change}})$ triplets. One such table was created for each possible (Δ, N) pair. Each value $p(R=\text{"change"} | N, \Delta, \bar{J}, \tau, p_{\text{change}})$ in these tables was approximated using the optimal decision rule (see Supplementary Information) applied to 10,000 Monte Carlo samples. Next, for a given set of parameter values, the log likelihood of each trial response was computed in two steps. First, the expected total loss was computed as a function of \bar{J} , using $\bar{L}_{\text{total}}(\bar{J}, N) = p_{\text{incorrect}}(\bar{J}, N) + \tilde{\lambda} \bar{J} N$, where $p_{\text{incorrect}}(\bar{J}, N)$ was estimated using the pre-computed tables. Second, we looked up $\log p(R_i | N_i, \Delta_i, \bar{J}_{\text{optimal}}, \tau, p_{\text{change}})$ from the pre-computed tables, where \bar{J}_{optimal} is the value of \bar{J} for which expected total loss was lowest. To estimate the best-fitting parameters, we performed a grid search over a large set of parameter combinations, separately for each subject.

Change localization and visual search. Model fitting methods for the change-localization and visual-search tasks were identical to the methods for the change-detection task, except for differences in the parameter vector (no prior in the change localization task; p_{present} instead of p_{change} in visual search) and the optimal decision rules (see Supplementary Information).

REFERENCES

1. Ma, W. J., Husain, M. & Bays, P. M. Changing concepts of working memory. *Nat. Neurosci.* **17**, 347–56 (2014).
2. Mazyar, H., van den Berg, R. & Ma, W. J. Does precision decrease with set size? *J. Vis.* **12**, 10 (2012).
3. Palmer, J. Attentional limits on the perception and memory of visual information. *J. Exp.*

- 321 *Psychol. Hum. Percept. Perform.* **16**, 332–350 (1990).
- 322 4. Shaw, M. L. in *Attention and performance VIII* (ed. Nickerson, R. S.) 277–296 (Erlbaum,
- 323 1980).
- 324 5. Lindsay, P. H., Taylor, M. M. & Forbes, S. M. Attention and multidimensional
- 325 discrimination. *Percept. Psychophys.* **4**, 113–117 (1968).
- 326 6. Zhang, W. & Luck, S. J. Discrete fixed-resolution representations in visual working
- 327 memory. *Nature* **453**, 233–235 (2008).
- 328 7. Sims, C. R., Jacobs, R. A. & Knill, D. C. An ideal observer analysis of visual working
- 329 memory. *Psychol. Rev.* **119**, 807–30 (2012).
- 330 8. Ma, W. J. & Huang, W. No capacity limit in attentional tracking: evidence for
- 331 probabilistic inference under a resource constraint. *J. Vis.* **9**, 3.1-30 (2009).
- 332 9. Bays, P. M. Noise in neural populations accounts for errors in working memory. *J.*
- 333 *Neurosci.* **34**, 3632–45 (2014).
- 334 10. Keshvari, S., van den Berg, R. & Ma, W. J. No Evidence for an Item Limit in Change
- 335 Detection. *PLoS Comput. Biol.* **9**, (2013).
- 336 11. van den Berg, R., Awh, E. & Ma, W. J. Factorial comparison of working memory models.
- 337 *Psychol. Rev.* **121**, 124–49 (2014).
- 338 12. Bays, P. M., Catalao, R. F. G. & Husain, M. The precision of visual working memory is
- 339 set by allocation of a shared resource. *J. Vis.* **9**, 7.1-11 (2009).
- 340 13. Bays, P. M. & Husain, M. Dynamic shifts of limited working memory resources in human
- 341 vision. *Science (80-.).* **321**, 851–4 (2008).
- 342 14. van den Berg, R., Shin, H., Chou, W.-C., George, R. & Ma, W. J. Variability in encoding
- 343 precision accounts for visual short-term memory limitations. *Proceedings of the National*
- 344 *Academy of Sciences* **109**, 8780–8785 (2012).
- 345 15. Devkar, D. T. & Wright, A. A. The same type of visual working memory limitations in
- 346 humans and monkeys. *J. Vis.* **13**, 1–18 (2015).
- 347 16. Elmore, L. C. *et al.* Visual short-term memory compared in rhesus monkeys and humans.
- 348 *Curr. Biol.* **21**, 975–979 (2011).
- 349 17. Wilken, P. & Ma, W. J. A detection theory account of change detection. *J. Vis.* **4**, 1120–35
- 350 (2004).
- 351 18. Oberauer, K., Farrell, S., Jarrold, C. & Lewandowsky, S. What Limits Working Memory

- Capacity? *Psychol. Bull.* **142**, 758–799 (2016).
19. Oberauer, K. & Lin, H. An Interference Model of Visual Working Memory. *Psychol. Rev.* **124**, 21–59 (2017).
20. Attwell, D. & Laughlin, S. B. An energy budget for signaling in the grey matter of the brain. *J. Cereb. Blood Flow Metab.* **21**, 1133–1145 (2001).
21. Lennie, P. The cost of cortical computation. *Curr. Biol.* **13**, 493–497 (2003).
22. Sterling, P. & Laughlin, S. *Principles of neural design*. (MIT Press, 2015).
23. Pestilli, F. & Carrasco, M. Attention enhances contrast sensitivity at cued and impairs it at uncued locations. *Vision Res.* **45**, 1867–1875 (2005).
24. Christie, S. T. & Schrater, P. Cognitive cost as dynamic allocation of energetic resources. *Front. Neurosci.* **9**, (2015).
25. Della Libera, C. & Chelazzi, L. Visual selective attention and the effects of monetary rewards. *Psychol. Sci. a J. Am. Psychol. Soc. / APS* **17**, 222–227 (2006).
26. Peck, C. J., Jangraw, D. C., Suzuki, M., Efem, R. & Gottlieb, J. Reward modulates attention independently of action value in posterior parietal cortex. *J. Neurosci.* **29**, 11182–11191 (2009).
27. Baldassi, S. & Simoncini, C. Reward sharpens orientation coding independently of attention. *Front. Neurosci.* (2011). doi:10.3389/fnins.2011.00013
28. Mazyar, H., Van den Berg, R., Seilheimer, R. L. & Ma, W. J. Independence is elusive : Set size effects on encoding precision in visual search. *J. Vis.* **13**, 1–14 (2013).
29. Keshvari, S., van den Berg, R. & Ma, W. J. Probabilistic computation in human perception under variability in encoding precision. *PLoS One* **7**, (2012).
30. Cover, T. M. & Thomas, J. A. *Elements of Information Theory. Elements of Information Theory* (2005). doi:10.1002/047174882X
31. Fougny, D., Suchow, J. W. & Alvarez, G. A. Variability in the quality of visual working memory. *Nat. Commun.* **3**, 1229 (2012).
32. Ma, W. J., Beck, J. M., Latham, P. E. & Pouget, A. Bayesian inference with probabilistic population codes. *Nat. Neurosci.* **9**, 1432–1438 (2006).
33. Paradiso, M. a. A theory for the use of visual orientation information which exploits the columnar structure of striate cortex. *Biol. Cybern.* **58**, 35–49 (1988).
34. Seung, H. S. & Sompolinsky, H. Simple models for reading neuronal population codes.

Proc.Natl.Acad.Sci. **90**, 10749–10753 (1993).

35. Mankiw, N. G. *Principles of economics. Book 328*, (2004).

36. Akaike, H. A new look at the statistical model identification. *IEEE Trans. Automat. Contr.* **19**, (1974).

37. Sims, C. R. The cost of misremembering: Inferring the loss function in visual working memory. *J. Vis.* **15**, 2 (2015).

38. Dyrholm, M., Kyllingsbæk, S., Espeseth, T. & Bundesen, C. Generalizing parametric models by introducing trial-by-trial parameter variability: The case of TVA. *J. Math. Psychol.* **55**, 416–429 (2011).

39. Wolpert, D. M. & Landy, M. S. Motor control is decision-making. *Current Opinion in Neurobiology* **22**, 996–1003 (2012).

40. Green, D. M. & Swets, J. A. Signal detection theory and psychophysics. *Society* **1**, 521 (1966).

41. Körding, K. Decision theory: what ‘should’ the nervous system do? *Science* **318**, 606–610 (2007).

42. Geisler, W. S. Contributions of ideal observer theory to vision research. *Vision Research* **51**, 771–781 (2011).

43. Shen, S. & Ma, W. J. A detailed comparison of optimality and simplicity in perceptual decision making. *Psychol. Rev.* **123**, 452–480 (2016).

44. Vincent, B. T., Baddeley, R. J., Troscianko, T. & Gilchrist, I. D. Is the early visual system optimised to be energy efficient? *Network* **16**, 175–190 (2005).

45. Liu, Y. S., Stevens, C. F. & Sharpee, T. Predictable irregularities in retinal receptive fields. *Proc. Natl. Acad. Sci.* **106**, 16499–16504 (2009).

46. Olshausen, B. A. & Field, D. J. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature* **381**, 607–609 (1996).

47. Attneave, F. Some informational aspects of visual perception. *Psychol. Rev.* **61**, 183–193 (1954).

48. Barlow, H. B. H. in *Sensory Communication* 217–234 (1961).
doi:10.1080/15459620490885644

49. Ganguli, D. & Simoncelli, E. P. Implicit encoding of prior probabilities in optimal neural populations. *Adv. Neural Inf. Process. Syst.* **2010**, 658–666 (2010).

50. Cherniak, C. Component placement optimization in the brain. *J. Neurosci.* **14**, 2418–2427 (1994).
51. Chklovskii, D. B., Schikorski, T. & Stevens, C. F. Wiring optimization in cortical circuits. *Neuron* **34**, 341–347 (2002).
52. Clune, J., Mouret, J.-B. & Lipson, H. The evolutionary origins of modularity. *Proc. R. Soc. B Biol. Sci.* **280**, 20122863–20122863 (2013).
53. Sims, C. R. Rate–distortion theory and human perception. *Cognition* **152**, 181–198 (2016).
54. Wei, Z., Wang, X.-J. & Wang, D.-H. From Distributed Resources to Limited Slots in Multiple-Item Working Memory: A Spiking Network Model with Normalization. *J. Neurosci.* **32**, 11228–11240 (2012).
55. Orhan, A. E. & Ma, W. J. Neural Population Coding of Multiple Stimuli. *J. Neurosci.* **35**, 3825–3841 (2015).
56. Endress, A. & Szabó, S. Interference and memory capacity limitations. *Psychol. Rev.* **In press**, (2017).
57. Simon, H. A. *Models of Man (Book)*. *Operations Research* **5**, (1957).
58. Krajbich, I. & Rangel, A. Multialternative drift-diffusion model predicts the relationship between visual fixations and choice in value-based decisions. *Proc. Natl. Acad. Sci.* **108**, 13852–13857 (2011).
59. Sims, C. A. Implications of rational inattention. *J. Monet. Econ.* **50**, 665–690 (2003).
60. Faisal, A. A., Selen, L. P. J. & Wolpert, D. M. Noise in the nervous system. *Nat. Rev. Neurosci.* **9**, 292–303 (2008).
61. Beck, J. M., Ma, W. J., Pitkow, X., Latham, P. E. & Pouget, A. Not Noisy, Just Wrong: The Role of Suboptimal Inference in Behavioral Variability. *Neuron* **74**, 30–39 (2012).
62. Yi, D.-J., Woodman, G. F., Widders, D., Marois, R. & Chun, M. M. Neural fate of ignored stimuli: dissociable effects of perceptual and working memory load. *Nat. Neurosci.* **7**, 992–996 (2004).
63. Bae, G., Allred, S. R., Wilson, C. & Flombaum, J. I. Stimulus-specific variability in color working memory with delayed estimation. *J. Vis.* **14**, 1–23 (2014).
64. Girshick, A. R., Landy, M. S. & Simoncelli, E. P. Cardinal rules: visual orientation perception reflects knowledge of environmental statistics. *Nat. Neurosci.* **14**, 926–932

(2011).

65. Wei, X.-X. & Stocker, A. A. A Bayesian observer model constrained by efficient coding can explain ‘anti-Bayesian’ percepts. *Nat. Neurosci.* **18**, 1509–1517 (2015).
66. Fuster, J. M. & Alexander, G. E. Neuron Activity Related to Short-Term Memory. *Science* **173**, 652–654 (1971).
67. Funahashi, S., Bruce, C. J. & Goldman-Rakic, P. S. Mnemonic coding of visual space in the monkey’s dorsolateral prefrontal cortex. *J. Neurophysiol.* **61**, 331–349 (1989).
68. Mongillo, G., Barak, O. & Tsodyks, M. Synaptic Theory of Working Memory. *Science* (80-.). **319**, 1543–1546 (2008).
69. Zhang, W. & Luck, S. J. The Number and Quality of Representations in Working Memory. *Psychol. Sci.* **22**, 1434–1441 (2011).
70. Carandini, M. & Heeger, D. Normalization as a canonical neural computation. *Nat. Rev. Neurosci.* 1–12 (2012). doi:10.1038/nrn3136
71. Reynolds, J. H. & Heeger, D. J. The Normalization Model of Attention. *Neuron* **61**, 168–185 (2009).