

1 **A rational theory of set size effects in working memory and attention**

2
3 Ronald van den Berg¹ and Wei Ji Ma²

4
5 ¹Department of Psychology, University of Uppsala, Uppsala, Sweden.

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

8
9 **Classification**

10 SOCIAL SCIENCES – Psychological and Cognitive Sciences

11
12 **Keywords**

13 Working memory, attention, ecological rationality, models

14
15 **Characters (with spaces):** ~39,000

16
17 **Corresponding author**

18 Ronald van den Berg

19 Department of Psychology

20 Von Kraemers Allé 1A, 75237

21 Uppsala, Sweden

22 Tel: +46184717277

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

24

25 **The precision with which items are encoded in working memory and attention decreases with**
26 **the number of encoded items. Current theories typically account for this “set size effect” by**
27 **postulating a hard constraint on the allocated amount of some kind of encoding resource,**
28 **such as samples, spikes, slots, or bits. While these theories have produced models that are**
29 **descriptively successful, they offer no principled explanation for the very existence of set size**
30 **effects: given their detrimental consequences for behavioral performance, why have these**
31 **effects not been weeded out by evolutionary pressure, for example by allocating resources**
32 **proportionally to the number of encoded items? Here, we propose a theory that is based on**
33 **an ecological notion of rationality: set size effects are the result of an optimal trade-off**
34 **between behavioral performance and the neural costs associated with stimulus encoding. We**
35 **derive models for four visual working memory and attention tasks and show that they**
36 **account well for data from eleven previously published experiments. Our results suggest that**
37 **set size effects have a rational basis and that ecological costs should be considered in models**
38 **of human behavior.**

39
40 Cognitive performance is strongly constrained by set size effects in working memory and attention:
41 the precision with which these systems encode information rapidly declines with the number of
42 items, as observed in for example delayed estimation, change detection, visual search, and
43 multiple-object tracking tasks (1–8). By contrast, set size effects seem to be absent in long-term
44 memory, where fidelity has been found to be independent of set size (9). The existence of set size
45 effects is thus not a general property of neural coding, but rather a phenomenon that requires
46 explanation. Despite an abundance of models, such an explanation is still lacking.

47 A common way to model set size effects has been to postulate that stimuli are encoded
48 using a fixed total amount of resources, formalized as “samples” (1, 3, 10), slots (11), information
49 bit rate (12), Fisher information (8), or neural firing (13): the larger the number of encoded items,
50 the lower the amount of resource available for each item and, therefore, the lower the precision
51 per item. There are two problems with this explanation. Firstly, they predict that precision is
52 inversely proportional to set size, which is often not the case (e.g., (14–16)). Secondly, it is unclear
53 whether keeping the amount of encoding resources constant across set sizes serves any
54 ecologically relevant function: why has the brain not evolved to counter set size effects by
55 increasing the amount of allocated resource as more items are encoded, as seems to be the case in

56 long-term memory? To address the first problem, models have been proposed in which encoding
57 precision is postulated to be a power-law function of set size (5, 6, 15–21). These models tend to
58 provide excellent fits, but have been criticized for lacking a principled motivation for the
59 postulated power-law (22, 23), thus failing to address the second problem.

60 Here, we take an ecological approach to explain set size effects, starting from the principle
61 that neural firing is energetically costly (24–26). This cost may have pressured the brain to balance
62 behavioral benefits of high precision against the neural loss that it induces (8, 25, 27, 28). What
63 level of encoding precision establishes a good balance might depend on multiple factors, such as
64 set size, task, and motivation. Indeed, performance on perceptual decision-making tasks can be
65 improved by increasing monetary reward (29–31), which suggests that the total amount of resource
66 spent on encoding has some flexibility that is driven by ecological factors. Based on these
67 considerations, we hypothesize that set size effects on encoding precision reflect an ecologically
68 rational strategy that balances behavioral performance against neural costs. Below, we derive
69 formal models from this hypothesis for four visual working memory and attention tasks, fit them
70 to data from eleven previously published experiments, and discuss implications of our findings.

71

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

Experiment	Reference	Task	Feature	Set sizes	#subj
DE1	(6)	Delayed estimation	Color	1, 2, 4, 8	15
DE2	(11)	Delayed estimation	Color	1, 2, 3, 6	8
DE3	(17)	Delayed estimation	Color	1, 2, 4, 6	12
DE4	(18)	Delayed estimation	Orientation	1-8	6
DE5	(18)	Delayed estimation	Color	1-8	13
DE6	(18)	Delayed estimation	Color	1-8	13
CD1	(14)	Change detection	Color	1, 2, 4, 8	7
CD2	(14)	Change detection	Orientation	2, 4, 6, 8	10
CL1	(18)	Change localization	Color	2, 4, 6, 8	7
CL2	(18)	Change localization	Orientation	2, 4, 6, 8	11
VS	(4)	Visual search	Orientation	1, 2, 4, 8	6

75

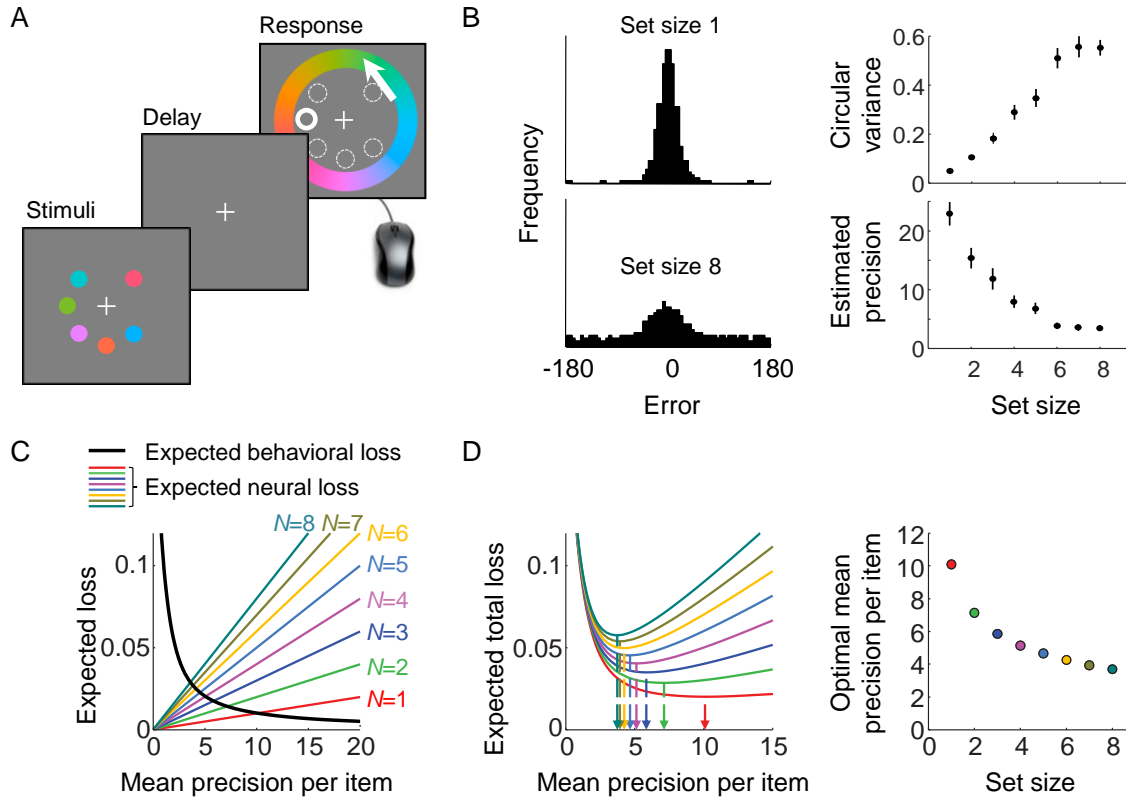


Figure 1. An ecologically rational model of set size effects in delayed estimation. (A) Example of a delayed-estimation experiment. The subject is briefly presented with a set of stimuli and, after a short delay, reports the value of a randomly chosen target item. (B) Estimation error distributions widen 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 $\tau \downarrow 0$.

76

77

78 THEORY

79 We first present our theory in the context of the delayed-estimation paradigm (6) and will later
 80 show how it generalizes to other tasks. In single-probe delayed-estimation tasks, subjects briefly
 81 hold a set of items in memory and report their estimate of a randomly chosen target item (Fig. 1A;
 82 Table 1). Estimation error ε is the (circular) difference between the subject's estimate and the true
 83 stimulus value s . Set size effects in this task are visible as a widening of the estimation error
 84 distribution (Fig. 1B). As in previous work (4, 5, 14, 15, 18, 32), we assume that a memory x
 85 follows a Von Mises distribution with mean s and concentration parameter κ , and define encoding

86 precision J as Fisher information (33), which is one-to-one related to κ (see Supplementary
87 Information). We assume that response noise is negligible, such that the estimation error is equal
88 to the memory error, $\varepsilon=x-s$. Moreover, we assume variability in J across items and trials (5, 15,
89 18, 32, 34), which we model using a gamma distribution with a mean \bar{J} and a scale parameter τ
90 (see Supplementary Information).

91 The key novelty of our theory is the idea that stimuli are encoded with a level of mean
92 precision, \bar{J} , that minimizes a combination of *behavioral loss* and *neural loss*. Behavioral loss is
93 induced by making an error ε , which we formalize using a mapping $L_{\text{behavioral}}(\varepsilon)$. This mapping may
94 depend on both internal incentives (e.g., intrinsic motivation) and external ones (e.g., the reward
95 scheme imposed by the experimenter). For the moment, we choose a power-law function,
96 $L_{\text{behavioral}}(\varepsilon)=|\varepsilon|^\beta$ with $\beta>0$ as a free parameter, such that larger errors correspond with larger loss.
97 The *expected* behavioral loss, denoted $\bar{L}_{\text{behavioral}}$, is obtained by averaging loss across all possible
98 errors, weighted by the probability that each error occurs,

99

$$100 \quad \bar{L}_{\text{behavioral}}(\bar{J}, N) = \int L_{\text{behavioral}}(\varepsilon) p(\varepsilon | \bar{J}, N) d\varepsilon, \quad (1)$$

101

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

105 The second kind of loss is a *neural loss* induced by the neural spiking activity that represents
106 a stimulus. For many choices of spike variability, including the common one of Poisson-like
107 variability (35), the precision (Fisher information) of a stimulus encoded in a neural population is
108 proportional to the neural spiking rate (36, 37). Moreover, it has been estimated that the energetic
109 loss induced by each spike increases with spiking rate (24, 25). When combining these two
110 premises, the expected neural loss associated with the encoding of an item is a supralinear function
111 of encoding precision, which can be modeled using for example a power-law function,
112 $\bar{L}_{\text{neural}}(\bar{J}) = \alpha \bar{J}^\omega$, with α and ω as free parameters. However, to minimize the number of free
113 parameters, we assume for the moment that the function is linear ($\beta=1$) and will later present a
114 mathematical proof that our theory generalizes to supralinear functions ($\beta>0$; condition (iv) at the

115 end of this section). Further assuming that stimuli are encoded independently of each other, neural
116 loss is also proportional to the number of encoded items, N . We thus obtain

117
118
$$\bar{L}_{\text{neural}}(\bar{J}, N) = \alpha \bar{J} N, \quad (2)$$

119
120 where α is a free parameter that represents the amount of neural loss incurred by a unit increase in
121 mean precision (Fig. 1C, colored lines).

122 We combine the two types of expected loss into a total expected loss function (Fig. 1D),

123
124
$$\begin{aligned} \bar{L}_{\text{total}}(\bar{J}, N) &= \bar{L}_{\text{behavioral}}(\bar{J}, N) + \lambda \bar{L}_{\text{neural}}(\bar{J}, N) \\ &= \bar{L}_{\text{behavioral}}(\bar{J}, N) + \lambda \alpha \bar{J} N, \end{aligned} \quad (3)$$

125
126 where the weight $\lambda \geq 0$ represents the importance of keeping neural loss low relative to the
127 importance of good performance. Since λ and α have interchangeable effects on the model
128 predictions, they can be fitted as a single free parameter $\tilde{\lambda} \triangleq \lambda \alpha$. We refer to the level of mean
129 precision that minimizes the total expected loss as *optimal mean precision*,

130
131
$$\bar{J}_{\text{optimal}}(\bar{J}, N) = \underset{J}{\operatorname{argmin}} \bar{L}_{\text{total}}(\bar{J}, N). \quad (4)$$

132
133 Under the loss functions proposed above, we find that \bar{J}_{optimal} is a decreasing function of set size
134 (Fig. 1D), which is qualitatively consistent with set size effects observed in experimental data (cf.
135 Fig. 1B). However, it can be shown (see Supplementary Information) that the conditions under
136 which this model predicts a set size effect generalize to any choice of loss functions, as long as the
137 four, rather general conditions are satisfied: (i) the expected behavioral loss is a strictly decreasing
138 function of encoding precision, i.e., an increase in precision results in an increase in performance;
139 (ii) the expected behavioral loss is subject to a law of diminishing returns (38): the higher the initial
140 precision, the smaller the behavioral benefit obtained from an increase in precision; (iii) the
141 expected neural loss is an increasing function of encoding precision; (iv) the expected neural loss
142 associated with a fixed increase in precision increases with precision. Next, we evaluate the model
143 by fitting it to data from a range of previously published experiments.

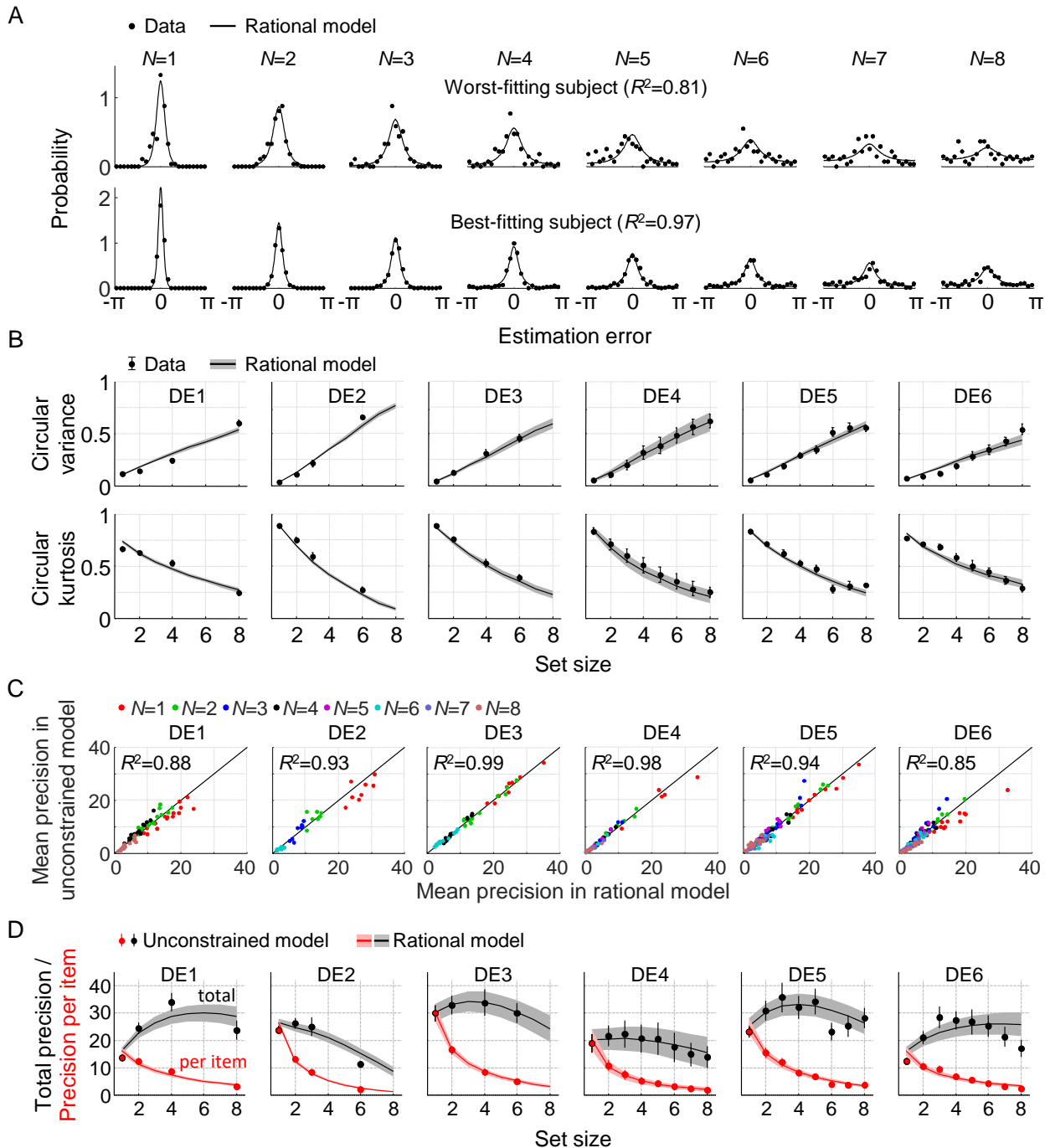


Figure 2. Model fits to data from six delayed-estimation experiments. (A) Maximum-likelihood fits to raw data of the worst-fitting and best-fitting subjects (based on R^2). (B) Subject-averaged fits to the two statistics that summarize the estimation error distributions (circular variance kurtosis) as 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.

145 **RESULTS**

146

147 **Model fits**

148 We used maximum-likelihood estimation to fit the three model parameters ($\tilde{\lambda}$, τ , and β) to 67
149 individual-subject data sets from a delayed-estimation benchmark set* (Table 1). The model
150 accounts well for the raw error distributions (Fig. 2A) and the two statistics that summarize these
151 distributions (Fig. 2B). Model comparison based on the Akaike Information Criterion (AIC) (39)
152 indicates that the goodness of fit is comparable to that of a descriptive model variant in which the
153 relation between encoding precision and set size is assumed to follow a power law
154 ($\Delta\text{AIC}=5.27\pm 0.70$ in favor of the rational model). Hence, the rational model provides a principled
155 explanation of set size effects in delayed-estimation tasks without sacrificing quality of fit.

156

157 **Comparison with an unconstrained model**

158 We next try to falsify our theory by testing whether a mapping between set size and encoding
159 precision can be found that fits the data better than the relation imposed by the loss-minimization
160 strategy of the rational model. To this end, we fit an unconstrained variant of the model in which
161 memory precision is fitted as a free parameter at each set size. We find only a minimal difference
162 in goodness of fit ($\Delta\text{AIC}=3.49\pm 0.93$ in favor of the unconstrained model), suggesting that the fits
163 of the rational model are close to the best possible fits. This finding is corroborated by examination
164 of the fitted parameter values: the estimated precision values in the unconstrained model closely
165 match the precision values in the rational model (Fig. 2C). Hence, it seems that no relation exists
166 that fits these data substantially better than the constrained set of relations that are possible in the
167 rational model.

168

169 **Total amount of allocated resource as a function of set size**

170 We estimate the total amount of allocated encoding resource as the mean precision (Fisher
171 information) per item summed across all items, $\bar{J}_{\text{total}} \triangleq \bar{J}N$. In fixed-resource models \bar{J}_{total} is by

* The original benchmark set (15) 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.

172 definition constant and in power-law models it is a monotonic function of set size. However, an
173 interesting qualitative feature of the rational model is that in some of the experiments the best-
174 fitting parameters produce a *non-monotonic* relation between \bar{J}_{total} and set size (Fig. 2D, gray
175 curves). This means that at small set sizes it apparently sometimes pays off (in terms of total-loss
176 minimization) to increase the total amount of allocated resource when an item is added, while the
177 opposite is true at large set sizes[†]. Using the best-fitting precision values in the unconstrained
178 model as an estimate of how much encoding resource subjects allocated on average at each set
179 size, we find that the data show clear signs of a similar non-monotonicity (Fig. 2D, black circles);
180 to our knowledge, this has not previously been reported.

181

182 **Alternative loss functions**

183 To evaluate the necessity of a free parameter in the behavioral loss function, $L_{\text{behavioral}}(\varepsilon)$, we also
184 test the following three parameter-free choices: $|\varepsilon|$, ε^2 , and $-\cos(\varepsilon)$. Model comparison favors the
185 original model with AIC differences of 14.0 ± 2.8 , 24.4 ± 4.1 , and 19.5 ± 3.5 , respectively. While there
186 may be other parameter-free functions that give better fits, we expect that a free parameter is
187 unavoidable here, as it is likely that the error-to-loss mapping differs across experiments (due to
188 differences in external incentives) and possibly also across subjects within an experiment (due to
189 differences in internal incentives). We also test a two-parameter function that was proposed
190 recently (Eq. (5) in (40)). The main difference with our original choice is that this alternative
191 function allows for saturation effects in the error-to-loss mapping. However, this extra flexibility
192 does not increase the goodness of fit sufficiently to justify the additional parameter, as the original
193 model outperforms this variant with an AIC difference of 5.3 ± 1.8 .

194

195 **Generalization to other tasks**

196 We next examine the generality of our theory, by testing whether it can also explain set size effects
197 in two change detection tasks (Table 1). In these experiments, the subject is on each trial
198 sequentially presented with two sets of stimuli and reports whether there was a change at any of
199 the stimulus locations (Fig. 3A). A change was present on half of the trials, at a random location

[†] Upon reflection, it is perhaps not surprising to occasionally find a non-monotonic relation: when multiplying a decreasing function of set size (\bar{J} as a function of N) with an increasing one (N itself), it is easy to obtain a function that is not monotonic but peaks at an intermediate value.

200 and with a random change magnitude. The behavioral error, ε , takes only two values in this task:
201 “correct” and “incorrect”. Therefore, $p(\varepsilon | \bar{J}, N)$ specifies the probabilities of correct and
202 incorrect responses for a given level of precision and set size, which depend on the observer’s
203 decision rule. Following previous work (14, 32), we assume that subjects use the Bayes-optimal
204 rule (see Supplementary Information) and that there is random variability in encoding precision.
205 This decision rule introduces one free parameter, p_{change} , specifying the subject's degree of prior
206 belief that a change will occur. Due to the binary nature of ε in this task, the free parameter of the
207 behavioral loss function drops out of the model, as its effect is equivalent to changing parameter
208 $\tilde{\lambda}$ (see Supplementary Information). The model thus has three free parameters ($\tilde{\lambda}$, τ , and p_{change}).
209 We find that the maximum-likelihood fits account well for the data in both experiments (Fig. 3B).

210 So far, we have considered tasks with continuous and binary judgments. We next consider
211 two change localization experiments (Table 1) in which judgments are non-binary but categorical.
212 The task is identical to change detection, except that a change is present on every trial and the
213 observer reports the location at which the change occurred (out of 2, 4, 6, or 8 locations). We again
214 assume variable precision and an optimal decision rule (see Supplementary Information).
215 Although the rational model has only two free parameters ($\tilde{\lambda}$ and τ), it accounts well for both
216 datasets (Fig. 3C).

217 The final task to which we apply our theory is a visual search experiment (4) (Table 1).
218 Unlike the previous three tasks, this is not a working memory task, as there was no delay period
219 between stimulus offset and response. Set size effects in this experiment are thus likely to stem
220 from limitations in attention rather than memory, but our theory applies without any additional
221 assumptions. Subjects judged whether a vertical target was present among one of N briefly
222 presented oriented ellipses (Fig. 3D). The distractors were drawn from a Von Mises distribution
223 centered at vertical. The width of the distractor distribution determined the level of heterogeneity
224 in the search display. Each subject was tested under three different levels of heterogeneity. We
225 again assume variable precision and an optimal decision rule (see Supplementary Information).
226 This decision rule has one free parameter, p_{present} , specifying the subject's prior degree of belief
227 that a target will be present. We fit the three free parameters ($\tilde{\lambda}$, τ , and p_{present}) to the data from all
228 three heterogeneity conditions at once and find that the model accounts well for the dependencies
229 of the hit and false alarm rates on both set size and distractor heterogeneity (Fig. 3E).

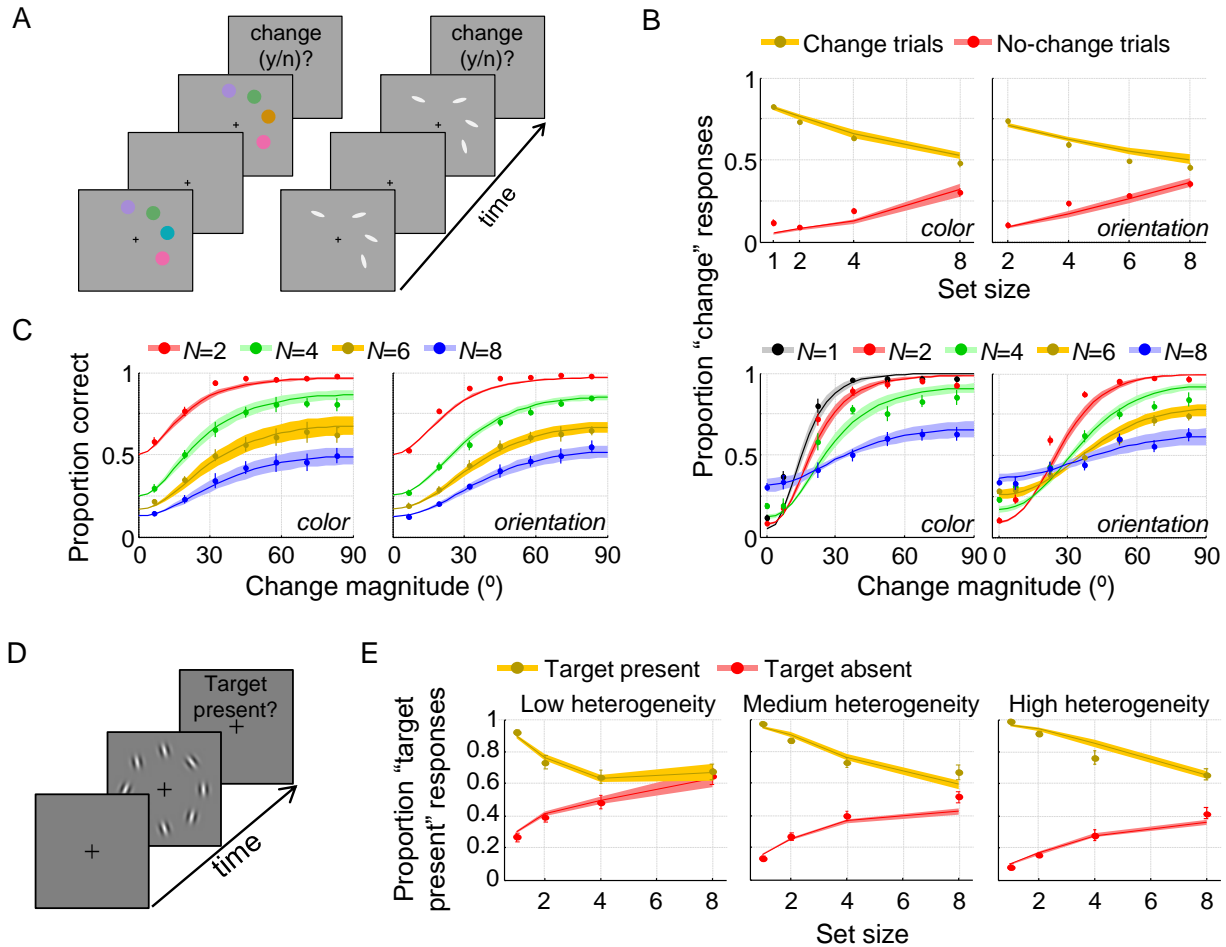


Figure 3. Model fits to three categorical decision-making tasks. (A) Experimental paradigm in the change-detection experiments. The paradigm for change localization was the same, except that a change was present on each trial and subjects reported the location of change. (B) Model fits to change-detection 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.

230

231

232 DISCUSSION

233 A key strength of our theory is that it uses a single principle of rationality and relatively few
 234 parameters to produce well-fitting models across a range of quite different tasks. Nevertheless,
 235 consideration of additional mechanisms could further improve the fits and lead to more complete
 236 models of human behavior. For example, previous studies have incorporated response noise (15,
 237 18), non-target responses (17), and a (variable) limit on the number of remembered items (12, 15,
 238 41) to improve fits. We did not consider such components here, as they come with additional
 239 parameters, some are task-specific (such as non-target responses), and they have so far not been

240 motivated in a principled manner. Regarding the latter point, it might be possible to treat some of
241 these mechanisms using an ecologically rational approach as well. For example, the level of
242 response noise might be set by optimizing a trade-off between performance and motor control
243 effort (42).

244 Our findings suggest that set size effects in working memory and attention may reflect a near-
245 optimal compromise reached by a system that strives to simultaneously maximize performance
246 and minimize spiking activity. A possible explanation why long-term memory does not seem to
247 suffer from set size effects (9) and has much larger capacity (43) is that loss incurred by
248 maintaining synaptic connections is likely to be lower than the loss incurred by persistent activity.

249 The work presented in this paper speaks to the relation between descriptive and rational
250 theories in psychology and neuroscience. The main motivation for rational theories is to reach a
251 deeper level of understanding by analyzing a system in the context of the ecological needs and
252 constraints that it evolved under. Besides the large literature on ideal-observer decision rules (44–
253 47), rational approaches have been used to explain properties of receptive fields (48–50), tuning
254 curves (51–53), neural wiring (54, 55), and neural network modularity (56). A transition from
255 descriptive to rational explanations might be an essential step in the maturation of theories of
256 biological systems, and in psychology there certainly seems to be more room for this kind of
257 explanation.

258 Although several previous models in the field of working memory and attention contain
259 rational aspects, none of them accounts for set size effects in a principled way. Sims and colleagues
260 have examined how errors in visual working memory can be minimized by optimally taking into
261 account statistics of the stimulus distribution, but assume a fixed total amount of available
262 encoding resource (12, 57). Moreover, in our own previous work on visual search (4, 5), change
263 detection (14, 32), and change localization (18), we used optimal-observer models for the decision
264 stage, but assumed an ad hoc power law for the encoding stage. An alternative explanation of set
265 size effects has been that the brain is unable to keep neural representations of multiple items
266 segregated from one another (23, 58–60): as the number of encoded items increases, so does the
267 level of interference in their representations, resulting in lower task performance. However, these
268 models offer no principled justification for the existence of interference and some require
269 additional mechanisms to account for set size effects; for example, the model by Oberauer and
270 colleagues requires three additional components – including a set-size dependent level of

271 background noise – to fully account for set size effects (23). That being said, our theory does not
272 rule out the possibility of interference, and it could be added onto any of the models we presented.

273 Our approach shares both similarities and differences with the concept of bounded rationality
274 (61), which states that human behavior is guided by mechanisms that provide “good enough”
275 solutions rather than optimal ones. The main similarity is that both approaches acknowledge that
276 human behavior is constrained by various cognitive limitations. However, an important difference
277 is that bounded rationality postulates these limitations as a given fact, while our approach explains
278 them as rational outcomes of ecological optimization processes. The suggestion that cognitive
279 limitations are themselves subject to optimization may also have implications for theories outside
280 the field of psychology. One example concerns recent models of value-based decision-making that
281 incorporate constraints imposed by working memory and attention limitations (e.g., (62)). Another
282 example is the theory of “rational inattention” in behavioral economics, which examines optimal
283 decision-making under the assumption that decision makers have a fixed limit on the total amount
284 of attention that they can allocate to process economic data (63). It might be interesting to extend
285 that theory by treating the amount of allocable attention as the outcome of an optimization process
286 rather than a constant.

287 While our results show that set size effects can in principle be explained as the result of an
288 optimization strategy, they do not necessarily imply that encoding precision is fully optimized on
289 every trial at any given task. First, encoding precision in the brain most likely has an upper limit,
290 due to irreducible sources of noise such as Johnson noise and Poisson shot noise (64), as well as
291 suboptimalities early in sensory processing (65). This prohibits subjects to reach the near-perfect
292 performance levels that our model may predict when the behavioral loss associated to errors is
293 huge. Second, precision might have a lower limit: task-irrelevant stimuli are sometimes
294 automatically encoded (66), perhaps because in natural environments few stimuli are ever
295 completely irrelevant. This would prevent subjects from sometimes encoding nothing at all, in
296 contradiction to what our theory predicts to happen at very large set sizes. Third, all models that
297 we tested incorporated variability in encoding precision. Part of this variability is possibly due to
298 stochastic factors such as neural noise, but part of it may also be systematic in nature (e.g.,
299 particular colors and orientations may be encoded with higher precision than others (67, 68)).
300 Whereas the systematic component could have a rational basis (e.g., higher precision for colors
301 and orientations that occur more frequently in natural scenes (53, 69)), this is unlikely to be true

302 for the random component. Indeed, when we jointly optimize \bar{J} and τ , we find estimates of τ that
303 are consistently 0, meaning that any variability in encoding precision is suboptimal from the
304 perspective of our model. Finally, even if set size effects are the result of a rational trade-off
305 between behavioral and neural loss, it may be that the solution that the brain settled on only works
306 well on average rather than being tailored to every possible situation. In that case, set size effects
307 could be more rigid across environmental changes (e.g., in task or reward structure) than predicted
308 by a fully optimal model that incorporates every such change.

309 Future work could further examine optimality of encoding precision in working memory and
310 attention by experimentally varying factors that affect the loss functions. In delayed estimation, an
311 obvious choice for this would be the delay period. Assuming that working memories are
312 maintained in persistent activity (70, 71), a longer delay would induce a higher cost and decrease
313 optimal encoding precision. Another experimental parameter that could be varied is the error-to-
314 loss mapping. A previous study that performed this manipulation found a behavioral effect in one
315 experiment, but did not vary set size (72). None of the experiments modeled here contained this
316 manipulation (DE4-DE6 imposed an explicit loss function but did not vary it; the other
317 experiments had no explicit scoring system). Future studies could measure effects of changes in
318 explicitly imposed scoring systems and test how well a rational model accounts for such effects.
319 Related to this, it would be relevant to examine whether subjects are able to internalize
320 experimental loss functions in the timespan of a single experiment and, if not, to further
321 characterize their "natural" loss functions (40). Another line of possible future work would be to
322 examine whether our theory can be generalized to the level of objects (73, 74).

323 Developmental work has shown that working memory capacity estimates change with age (75,
324 76). Viewed from the perspective of our proposed theory, this raises the question why the optimal
325 trade-off between behavioral and neural loss would change with age. A speculative answer could
326 be that a subject's encoding efficiency (formalized by parameter α in Eq. (2)) may improve during
327 childhood. An increase in encoding efficiency (i.e., lower α) has the same effect in our model as a
328 decrease in the set size (i.e., higher N), which we know is accompanied by an increase in encoding
329 precision. Hence, our model would predict subjects to increase encoding precision over time,
330 which is qualitatively consistent with the findings of the developmental studies.

331 Finally, our results raise the question what neural mechanisms could implement the kind of
332 near-optimal resource allocation strategy that is the core of our theory. Some form of divisive

333 normalization (13, 77) would be a likely candidate, as it has the effect of lowering the gain when
334 set size is larger. Moreover, divisive normalization is already a key operation in neural models of
335 attention (78) and visual working memory (13, 58).

336

337 **METHODS**

338 **Data and code sharing**

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

340

341 **Model fitting**

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

343 $\boldsymbol{\theta} = \{\tilde{\lambda}, \beta, \tau\}$ that maximizes the log likelihood function, $\sum_{i=1}^n \log p(\varepsilon_i | N_i, \boldsymbol{\theta})$, where n is the

344 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

345 on that trial. To reduce the risk of converging into a local maximum, initial parameter estimates

346 were chosen based on a coarse grid search over a large range of parameter values. The predicted

347 estimation error distribution for a given parameter vector $\boldsymbol{\theta}$ was computed as follows. First, \bar{J}_{optimal}

348 was computed by applying Matlab's `fminsearch` function to Eq. (5). In this process, the integrals

349 over ε and J were approximated numerically by discretizing the distributions of these variables

350 into 100 and 20 equal-probability bins, respectively. Next, the gamma distribution over precision

351 with mean \bar{J}_{optimal} and scale parameter τ was discretized into 20 equal-probability bins. Thereafter,

352 the predicted estimation error distribution was computed under the central value of each bin.

353 Finally, these 20 predicted distributions were averaged. We verified that our results are robust

354 under changes in the number of bins used in the numerical approximations.

355 *Change detection.* Model fitting in the change detection task consisted of finding parameter

356 vector $\boldsymbol{\theta} = \{\tilde{\lambda}, \tau, p_{\text{change}}\}$ that maximizes $\sum_{i=1}^n \log p(R_i | \Delta_i, N_i, \boldsymbol{\theta})$, where n is the number of trials in

357 the subject's data set, R_i is the response ("change" or "no change"), Δ_i the magnitude of change,

358 and N_i the set size on the i^{th} trial. For computational convenience, Δ was discretized into 30 equally

359 spaced bins. To find the maximum-likelihood parameters, we first created a table with predicted

360 probabilities of "change" responses for a large range of $(\bar{J}, \tau, p_{\text{change}})$ triplets. One such table was

361 created for each possible (Δ, N) pair. Each value $p(R=\text{“change”} \mid N, \Delta, \bar{J}, \tau, p_{\text{change}})$ in these tables
362 was approximated using the optimal decision rule (see Supplementary Information) applied to
363 10,000 Monte Carlo samples. Next, for a given set of parameter values, the log likelihood of each
364 trial response was computed in two steps. First, the expected total loss was computed as a function
365 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
366 pre-computed tables. Second, we looked up $\log p(R_i \mid N_i, \Delta_i, \bar{J}_{\text{optimal}}, \tau, p_{\text{change}})$ from the pre-
367 computed tables, where \bar{J}_{optimal} is the value of \bar{J} for which expected total loss was lowest. To
368 estimate the best-fitting parameters, we performed a grid search over a large set of parameter
369 combinations, separately for each subject.

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

374

375 REFERENCES

- 376 1. Shaw ML (1980) Identifying attentional and decision-making components in information
377 processing. *Attention and Performance VIII*, ed Nickerson RS (Erlbaum, Hillsdale, NJ,
378 NJ), pp 277–296.
- 379 2. Ma WJ, Husain M, Bays PM (2014) Changing concepts of working memory. *Nat*
380 *Neurosci* 17(3):347–56.
- 381 3. Palmer J (1990) Attentional limits on the perception and memory of visual information. *J*
382 *Exp Psychol Hum Percept Perform* 16(2):332–350.
- 383 4. Mazyar H, Van den Berg R, Seilheimer RL, Ma WJ (2013) Independence is elusive : Set
384 size effects on encoding precision in visual search. *J Vis* 13(5):1–14.
- 385 5. Mazyar H, van den Berg R, Ma WJ (2012) Does precision decrease with set size? *J Vis*
386 12(6):10.
- 387 6. Wilken P, Ma WJ (2004) A detection theory account of change detection. *J Vis*
388 4(12):1120–35.
- 389 7. Palmer J (1994) Set-size effects in visual search: The effect of attention is independent of
390 the stimulus for simple tasks. *Vision Res* 34(13). doi:10.1016/0042-6989(94)90128-7.

- 391 8. Ma WJ, Huang W (2009) No capacity limit in attentional tracking: evidence for
392 probabilistic inference under a resource constraint. *J Vis* 9(11):3.1-30.
- 393 9. Brady TF, Konkle T, Gill J, Oliva A, Alvarez G a (2013) Visual long-term memory has
394 the same limit on fidelity as visual working memory. *Psychol Sci* 24(6):981–90.
- 395 10. Lindsay PH, Taylor MM, Forbes SM (1968) Attention and multidimensional
396 discrimination. *Percept Psychophys* 4(2):113–117.
- 397 11. Zhang W, Luck SJ (2008) Discrete fixed-resolution representations in visual working
398 memory. *Nature* 453(7192):233–235.
- 399 12. Sims CR, Jacobs RA, Knill DC (2012) An ideal observer analysis of visual working
400 memory. *Psychol Rev* 119(4):807–30.
- 401 13. Bays PM (2014) Noise in neural populations accounts for errors in working memory. *J*
402 *Neurosci* 34(10):3632–45.
- 403 14. Keshvari S, van den Berg R, Ma WJ (2013) No Evidence for an Item Limit in Change
404 Detection. *PLoS Comput Biol* 9(2).
- 405 15. van den Berg R, Awh E, Ma WJ (2014) Factorial comparison of working memory models.
406 *Psychol Rev* 121(1):124–49.
- 407 16. Bays PM, Husain M (2008) Dynamic shifts of limited working memory resources in
408 human vision. *Science (80-)* 321(5890):851–4.
- 409 17. Bays PM, Catalao RFG, Husain M (2009) The precision of visual working memory is set
410 by allocation of a shared resource. *J Vis* 9(10):7.1-11.
- 411 18. van den Berg R, Shin H, Chou W-C, George R, Ma WJ (2012) Variability in encoding
412 precision accounts for visual short-term memory limitations. *Proc Natl Acad Sci*
413 109:8780–8785.
- 414 19. Devkar DT, Wright AA (2015) The same type of visual working memory limitations in
415 humans and monkeys. *J Vis* 13(2015):1–18.
- 416 20. Elmore LC, et al. (2011) Visual short-term memory compared in rhesus monkeys and
417 humans. *Curr Biol* 21(11):975–979.
- 418 21. Donkin C, Kary A, Tahir F, Taylor R (2016) Resources masquerading as slots: Flexible
419 allocation of visual working memory. *Cogn Psychol* 85:30–42.
- 420 22. Oberauer K, Farrell S, Jarrold C, Lewandowsky S (2016) What Limits Working Memory
421 Capacity? *Psychol Bull* 142(March):758–799.

- 422 23. Oberauer K, Lin H (2017) An Interference Model of Visual Working Memory. *Psychol*
423 *Rev* 124(1):21–59.
- 424 24. Attwell D, Laughlin SB (2001) An energy budget for signaling in the grey matter of the
425 brain. *J Cereb Blood Flow Metab* 21(10):1133–1145.
- 426 25. Lennie P (2003) The cost of cortical computation. *Curr Biol* 13(6):493–497.
- 427 26. Sterling P, Laughlin S (2015) *Principles of neural design*. (MIT Press).
- 428 27. Pestilli F, Carrasco M (2005) Attention enhances contrast sensitivity at cued and impairs it
429 at uncued locations. *Vision Res* 45(14):1867–1875.
- 430 28. Christie ST, Schrater P (2015) Cognitive cost as dynamic allocation of energetic
431 resources. *Front Neurosci* 9(JUL). doi:10.3389/fnins.2015.00289.
- 432 29. Della Libera C, Chelazzi L (2006) Visual selective attention and the effects of monetary
433 rewards. *Psychol Sci a J Am Psychol Soc / APS* 17(3):222–227.
- 434 30. Peck CJ, Jangraw DC, Suzuki M, Efem R, Gottlieb J (2009) Reward modulates attention
435 independently of action value in posterior parietal cortex. *J Neurosci* 29(36):11182–
436 11191.
- 437 31. Baldassi S, Simoncini C (2011) Reward sharpens orientation coding independently of
438 attention. *Front Neurosci* (FEB). doi:10.3389/fnins.2011.00013.
- 439 32. Keshvari S, van den Berg R, Ma WJ (2012) Probabilistic computation in human
440 perception under variability in encoding precision. *PLoS One* 7(6).
- 441 33. Cover TM, Thomas JA (2005) *Elements of Information Theory* doi:10.1002/047174882X.
- 442 34. Fougnie D, Suchow JW, Alvarez GA (2012) Variability in the quality of visual working
443 memory. *Nat Commun* 3:1229.
- 444 35. Ma WJ, Beck JM, Latham PE, Pouget A (2006) Bayesian inference with probabilistic
445 population codes. *Nat Neurosci* 9(11):1432–1438.
- 446 36. Paradiso M a (1988) A theory for the use of visual orientation information which exploits
447 the columnar structure of striate cortex. *Biol Cybern* 58(1):35–49.
- 448 37. Seung HS, Sompolinsky H (1993) Simple models for reading neuronal population codes.
449 *ProcNatAcadSci* 90(22):10749–10753.
- 450 38. Mankiw NG (2004) *Principles of economics* doi:10.1017/CBO9780511511455.
- 451 39. Akaike H (1974) A new look at the statistical model identification. *IEEE Trans Automat*
452 *Contr* 19(6). doi:10.1109/TAC.1974.1100705.

- 453 40. Sims CR (2015) The cost of misremembering: Inferring the loss function in visual
454 working memory. *J Vis* 15(3):2.
- 455 41. Dyrholm M, Kyllingsbæk S, Espeseth T, Bundesen C (2011) Generalizing parametric
456 models by introducing trial-by-trial parameter variability: The case of TVA. *J Math*
457 *Psychol* 55(6):416–429.
- 458 42. Wolpert DM, Landy MS (2012) Motor control is decision-making. *Curr Opin Neurobiol*
459 22(6):996–1003.
- 460 43. Brady TF, Konkle T, Alvarez GA, Oliva A (2008) Visual long-term memory has a
461 massive storage capacity for object details. *Proc Natl Acad Sci* 105(38):14325–14329.
- 462 44. Green DM, Swets JA (1966) Signal detection theory and psychophysics. *Society* 1:521.
- 463 45. Körding K (2007) Decision theory: what “should” the nervous system do? *Science*
464 318:606–610.
- 465 46. Geisler WS (2011) Contributions of ideal observer theory to vision research. *Vision Res*
466 51(7):771–781.
- 467 47. Shen S, Ma WJ (2016) A detailed comparison of optimality and simplicity in perceptual
468 decision making. *Psychol Rev* 123(4):452–480.
- 469 48. Vincent BT, Baddeley RJ, Troscianko T, Gilchrist ID (2005) Is the early visual system
470 optimised to be energy efficient? *Network* 16:175–190.
- 471 49. Liu YS, Stevens CF, Sharpee T (2009) Predictable irregularities in retinal receptive fields.
472 *Proc Natl Acad Sci* 106(38):16499–16504.
- 473 50. Olshausen BA, Field DJ (1996) Emergence of simple-cell receptive field properties by
474 learning a sparse code for natural images. *Nature* 381(6583):607–609.
- 475 51. Attneave F (1954) Some informational aspects of visual perception. *Psychol Rev*
476 61(3):183–193.
- 477 52. Barlow HBH (1961) Possible principles underlying the transformation of sensory
478 messages. *Sensory Communication*, pp 217–234.
- 479 53. Ganguli D, Simoncelli EP (2010) Implicit encoding of prior probabilities in optimal neural
480 populations. *Adv Neural Inf Process Syst* 2010(December 2010):658–666.
- 481 54. Cherniak C (1994) Component placement optimization in the brain. *J Neurosci*
482 14(April):2418–2427.
- 483 55. Chklovskii DB, Schikorski T, Stevens CF (2002) Wiring optimization in cortical circuits.

- 484 *Neuron* 34(3):341–347.
- 485 56. Clune J, Mouret J-B, Lipson H (2013) The evolutionary origins of modularity. *Proc R Soc*
486 *B Biol Sci* 280(1755):20122863–20122863.
- 487 57. Sims CR (2016) Rate–distortion theory and human perception. *Cognition* 152:181–198.
- 488 58. Wei Z, Wang X-J, Wang D-H (2012) From Distributed Resources to Limited Slots in
489 Multiple-Item Working Memory: A Spiking Network Model with Normalization. *J*
490 *Neurosci* 32(33):11228–11240.
- 491 59. Orhan AE, Ma WJ (2015) Neural Population Coding of Multiple Stimuli. *J Neurosci*
492 35(9):3825–3841.
- 493 60. Endress A, Szabó S (2017) Interference and memory capacity limitations. *Psychol Rev In*
494 press.
- 495 61. Simon HA (1957) *Models of Man (Book)* doi:10.2307/2281884.
- 496 62. Krajbich I, Rangel A (2011) Multialternative drift-diffusion model predicts the
497 relationship between visual fixations and choice in value-based decisions. *Proc Natl Acad*
498 *Sci* 108(33):13852–13857.
- 499 63. Sims CA (2003) Implications of rational inattention. *J Monet Econ* 50(3):665–690.
- 500 64. Faisal AA, Selen LPJ, Wolpert DM (2008) Noise in the nervous system. *Nat Rev Neurosci*
501 9:292–303.
- 502 65. Beck JM, Ma WJ, Pitkow X, Latham PE, Pouget A (2012) Not Noisy, Just Wrong: The
503 Role of Suboptimal Inference in Behavioral Variability. *Neuron* 74(1):30–39.
- 504 66. Yi D-J, Woodman GF, Widders D, Marois R, Chun MM (2004) Neural fate of ignored
505 stimuli: dissociable effects of perceptual and working memory load. *Nat Neurosci*
506 7(9):992–996.
- 507 67. Bae G, Allred SR, Wilson C, Flombaum JI (2014) Stimulus-specific variability in color
508 working memory with delayed estimation. *J Vis* 14(4):1–23.
- 509 68. Girshick AR, Landy MS, Simoncelli EP (2011) Cardinal rules: visual orientation
510 perception reflects knowledge of environmental statistics. *Nat Neurosci* 14(7):926–932.
- 511 69. Wei X-X, Stocker AA (2015) A Bayesian observer model constrained by efficient coding
512 can explain “anti-Bayesian” percepts. *Nat Neurosci* 18(10):1509–1517.
- 513 70. Fuster JM, Alexander GE (1971) Neuron Activity Related to Short-Term Memory.
514 *Science (80-)* 173(3997):652–654.

- 515 71. Funahashi S, Bruce CJ, Goldman-Rakic PS (1989) Mnemonic coding of visual space in
516 the monkey's dorsolateral prefrontal cortex. *J Neurophysiol* 61(2):331–349.
- 517 72. Zhang W, Luck SJ (2011) The Number and Quality of Representations in Working
518 Memory. *Psychol Sci* 22(11):1434–1441.
- 519 73. Luck SJ, Vogel EK (1997) The capacity of visual working memory for features and
520 conjunctions. *Nature* 390(6657):279–281.
- 521 74. Wheeler ME, Treisman AM (2002) Binding in short-term visual memory. *J Exp Psychol*
522 *Gen* 131(1):48–64.
- 523 75. Simmering VR, Perone S (2012) Working memory capacity as a dynamic process. *Front*
524 *Psychol* 3(January):567.
- 525 76. Simmering VR (2012) The development of visual working memory capacity during early
526 childhood. *J Exp Child Psychol* 111(4):695–707.
- 527 77. Carandini M, Heeger D (2012) Normalization as a canonical neural computation. *Nat Rev*
528 *Neurosci* (November):1–12.
- 529 78. Reynolds JH, Heeger DJ (2009) The Normalization Model of Attention. *Neuron*
530 61(2):168–185.
- 531
- 532
- 533
- 534