

1 Efficient sampling 2 and noisy decisions

3 Joseph Heng¹, Michael Woodford², Rafael Polania^{1*}

*For correspondence:

rafael.polania@hest.ethz.ch (RP)

4 ¹Decision Neuroscience Lab, Dept. of Health Sciences and Technology, ETH Zurich,
5 Switzerland; ²Department of Economics, Columbia University, New York, USA

7 **Abstract** The precision of human decisions is limited by both processing noise and basing
8 decisions on finite information. But what determines the degree of such imprecision? Here we
9 develop an efficient coding framework for higher-level cognitive processes, in which information
10 is represented by a finite number of discrete samples. We characterize the sampling process that
11 maximizes perceptual accuracy or fitness under the often-adopted assumption that full
12 adaptation to an environmental distribution is possible, and show how the optimal process
13 differs when detailed information about the current contextual distribution is costly. We tested
14 this theory on a numerosity discrimination task, and found that humans efficiently adapt to
15 contextual distributions, but in the way predicted by the model in which people must economize
16 on environmental information. Thus, understanding decision behavior requires that we account
17 for biological restrictions on information coding, challenging the often-adopted assumption of
18 precise prior knowledge in higher-level decision systems.

20 Introduction

21 It is well-established that sensory perception is imprecise, and moreover that the precision of com-
22 parative judgments regarding sensory magnitudes are not uniform over the domain of possible
23 stimuli. Increasing evidence suggests that observed patterns of non-uniformity in discrimination
24 thresholds can often be explained as reflecting a principle of efficient coding: the idea that informa-
25 tion is encoded in ways that minimize the costs of inaccurate decisions given biological constraints
26 on information acquisition (*Niven and Laughlin, 2008; Sharpee et al., 2014*). While early applica-
27 tions of efficient coding theory have primarily been to early stages of sensory processing (*Laughlin,*
28 *1981; Ganguli and Simoncelli, 2014; Wei and Stocker, 2015*), it is worth considering whether similar
29 principles may also shape the structure of internal representations of higher-level concepts, such
30 as the perceptions of value that underlie economic decision making (*Louie and Glimcher, 2012;*
31 *Polanía et al., 2019; Rustichini et al., 2017*). In this work, we contribute to the efficient coding
32 framework applied to cognition and behavior in several respects.

33 A first aspect concerns the range of possible internal representation schemes that should be
34 considered feasible, which determines the way in which greater precision of discrimination in one
35 part of the stimulus space requires less precision of discrimination elsewhere. Implementational
36 architectures proposed by *Ganguli and Simoncelli (2014)* or *Wei and Stocker (2015)* assume a pop-
37 ulation coding scheme in which different neurons have distinct "preferred" stimuli. While this is
38 clearly relevant for some kinds of low-level sensory features such as orientation, it is not obvi-
39 ous that this kind of internal representation is used in representing higher-level concepts such as
40 economic values. We instead develop an efficient coding theory for a case in which an extensive
41 magnitude (something that can be described by a larger or smaller number) is represented by the
42 total number of processing units, from among a population of units operating in parallel, that "vote"

43 in favor of the magnitude's being larger rather than small. The internal representation therefore
44 necessarily consists of a finite collection of binary signals.

45 Our restriction to representations made up of binary signals is in conformity with the obser-
46 vation that neural systems at many levels appear to transmit information via discrete stochastic
47 events (*Schreiber et al., 2002; Sharpee, 2017*). Moreover, cognitive models with this general struc-
48 ture have been argued to be relevant for higher-order decision problems such as value-based
49 choice. For example, it has been suggested that the perceived values of choice options are con-
50 structed by acquiring samples of evidence from memory regarding the emotions evoked by the
51 presented items (*Shadlen and Shohamy, 2016*). Related accounts suggest that when a choice must
52 be made between alternative options, information is acquired via discrete samples of information
53 that can be represented as binary responses (e.g., "yes/no" responses to queries) (*Norman, 1968;*
54 *Weber and Johnson, 2009*). The seminal *decision-by-sampling* (DbS) theory (*Stewart et al., 2006*) sim-
55 ilarly posits an internal representation of magnitudes relevant to a decision problem by tallies of
56 the outcomes of a set of binary comparisons between the current magnitude and alternative values
57 sampled from memory. The architecture that we assume for imprecise internal representations
58 has the general structure of proposals of these kinds; but we go beyond the above-mentioned in-
59 vestigations, in analyzing what an efficient coding scheme consistent with our general architecture
60 would be like.

61 A second aspect concerns the conclusions about efficient coding depending on the objective
62 for which the encoding system is assumed to be optimized. Information maximization theories
63 (*Laughlin, 1981; Ganguli and Simoncelli, 2014; Wei and Stocker, 2015*) assume that the objective
64 should be maximal mutual information between the true stimulus magnitude and the internal
65 representation. While this may be a reasonable assumption in the case of early sensory processing,
66 it is less obvious in the case of circuits involved more directly in decision making, and in the latter
67 case an obvious alternative is to ask what kind of encoding scheme will best serve to allow accurate
68 decisions to be made. In the theory that we develop here, our primary concern is with encoding
69 schemes that maximize a subject's probability of giving a correct response to a binary decision.
70 However, we compare the coding rule that would be optimal from this standpoint to one that
71 would maximize mutual information, or to one that would maximize the expected value of the
72 chosen item rather than the probability of choosing the larger item.

73 Third, and most importantly, we extend our theory of efficient coding to consider not merely the
74 nature of an efficient coding system for a single environmental frequency distribution assumed to
75 be permanently relevant – so that there has been ample time for the encoding rule to be optimally
76 adapted to that distribution of stimulus magnitudes – but also an efficient approach to adjusting
77 the encoding as the environmental frequency distribution changes. Prior discussions of efficient
78 coding have often considered the optimal choice of an encoding rule for a single environmental
79 frequency distribution, and derived quantitative predictions for an empirical stimulus distribution
80 that is assumed to represent a permanent feature of the natural environment (*Laughlin, 1981;*
81 *Ganguli and Simoncelli, 2014*). Such an approach may make sense for a theory of neural coding in
82 cortical regions involved in early-stage processing of sensory stimuli, but is less obviously appro-
83 priate for a theory of the processing of higher-level concepts such as economic value, where the
84 idea that there is a single, permanently relevant frequency distribution of magnitudes that may be
85 encountered is also much more doubtful. Moreover, recent evidence suggests that the encoding
86 of economically relevant magnitudes (the size of monetary gains or losses, the probability of dif-
87 ferent possible outcomes from a gamble, the length of the delay until a payment will be received)
88 varies with changes in the frequency distribution of such quantities used in a particular experiment
89 (*Stewart et al., 2015*).

90 Hence we treat the collection of information about the currently relevant environmental fre-
91 quency distribution, in order to appropriately adjust the encoding rule for new stimuli, and not
92 simply the collection of information about the current individual stimulus, as part of the computa-
93 tions that must be undertaken by the processing units that we model. Our definition of an efficient

94 coding scheme posits that it is desirable to economize both on the amount of information about
95 the current contextual distribution that is used, as well as on the amount of information about
96 the current stimulus that is used. We model both types of information collection as sampling pro-
97 cesses: the information used about the contextual distribution is provided by a finite sample from
98 that distribution, while the information used about the current stimulus is provided by a finite
99 sample of binary responses to queries about that stimulus. We can then quantify the two types of
100 resource constraint with which we are concerned as limits on the number of samples used of each
101 of these types.

102 In the case in which we assume no effective bound on the number of samples from the contex-
103 tual distribution that can be used in encoding the magnitude of a new stimulus, our theory reduces
104 to the kind of efficient coding theory considered by previous authors: for each possible contextual
105 frequency distribution, new stimuli are encoded using a rule that is efficient (in the sense of max-
106 imizing response accuracy, subject to a finite bound on the number of binary signals used to rep-
107 resent an individual magnitude) for that particular distribution. And as we show, the predictions
108 of our theory in this case are similar (at least in the limiting case of a large though finite number of
109 binary signals) to those of the efficient coding theories proposed in previous work (*Laughlin, 1981*;
110 *Ganguli and Simoncelli, 2014*; *Wei and Stocker, 2015*). If, instead, it is important to economize on
111 the number of samples from the contextual distribution used to encode each new stimulus, we
112 obtain a different result. In particular, we demonstrate that when this second consideration has a
113 great enough weight, the optimal encoding rule corresponds to the DbS algorithm of *Stewart et al.*
114 (*2006*).

115 A second goal of our work is to test the relevance of these different possible models of efficient
116 coding in the case of numerosity discrimination. Judgments of the comparative numerosity of two
117 visual displays provide a test case of particular interest given our objectives. On the one hand, a
118 long literature has argued that imprecision in numerosity judgments has a similar structure to psy-
119 chophysical phenomena in many low-level sensory domains (*Nieder and Dehaene, 2009*; *Nieder*
120 *and Miller, 2003*). This makes it reasonable to ask whether efficient coding principles may also
121 be relevant in this domain. At the same time, numerosity is plainly a more abstract feature of vi-
122 sual arrays than low-level properties such as local luminosity, contrast, or orientation, and so can
123 be computed only at a later stage of processing. Moreover, processing of numerical magnitudes
124 is a crucial element of many higher-level cognitive processes, such as economic decision making;
125 and it is arguable that many rapid or intuitive judgments about numerical quantities, even when
126 numbers are presented symbolically, are based on an “approximate number system” of the same
127 kind as is used in judgments of the numerosity of visual displays (*Piazza et al., 2007*; *Nieder and*
128 *Dehaene, 2009*). It has further been argued that imprecision in the internal representation of nu-
129 merical magnitudes may underly imprecision and biases in economic decisions (*Khaw et al., In*
130 *Press*; *Woodford, In Press*).

131 It is well-known that the precision of discrimination between nearby numbers of items de-
132 creases in the case of larger numerosities, in approximately the way predicted by *Weber's Law*,
133 and this is often argued to support of a model of imprecise coding based on a logarithmic transfor-
134 mation of the true number (*Nieder and Dehaene, 2009*; *Nieder and Miller, 2003*). However, while
135 the precision of internal representations of numerical magnitudes is arguably of great evolution-
136 ary relevance (*Butterworth et al., 2018*; *Nieder, 2020*), it is unclear why a specifically logarithmic
137 transformation of number information should be of adaptive value, and also whether the same
138 transformation is used independent of context (*Pardo-Vazquez et al., 2019*; *Brus et al., 2019*). We
139 report new experimental data on numerosity discrimination by human subjects, in the case of two
140 different frequency distributions for the numerosity of the presented stimuli, and show both that
141 the observed variation in discriminability over the stimulus range differs somewhat from the pre-
142 dictions of a logarithmic coding model, and that it changes when the distribution of stimuli used
143 in the experiment is different. We also compare the observed pattern of variation in discriminabil-
144 ity with the predictions of our efficient coding theory, under a variety of assumptions about both

145 the performance measure and the weight assigned to economizing on the number of samples re-
146 quired from the contextual distribution. We find that our data are most consistent with the DbS
147 model, which is to say, to the predictions of an efficient coding theory for which the performance
148 measure is the frequency of correct comparative judgments, and a substantial weight is placed on
149 reducing the required number of samples from the contextual distribution.

150 Results

151 A general efficient sampling framework

152 We consider a situation in which the objective magnitude of a stimulus with respect to some feature
153 can be represented by a quantity v . When the stimulus is presented to an observer, it gives rise to
154 an imprecise representation r in the nervous system, on the basis of which the observer produces
155 any required response. The internal representation r can be stochastic, with given values being
156 produced with conditional probabilities $p(r|v)$ that depend on the true magnitude. Here, we are
157 more specifically concerned with discrimination experiments, in which two stimulus magnitudes
158 v_1 and v_2 are presented, and the subject must choose which of the two is greater. We suppose
159 that each magnitude v_i has an internal representation r_i , drawn independently from a distribution
160 $p(r_i|v_i)$ that depends only on the true magnitude of that individual stimulus. The observer's choice
161 must be based on a comparison of r_1 with r_2 .

162 One way in which the cognitive resources recruited to make accurate discriminations may be
163 limited is in the variety of distinct internal representations that are possible. When the complexity
164 of feasible internal representations is limited, there will necessarily be errors in the identification of
165 the greater stimulus magnitude in some cases, even assuming an optimal *decoding rule* for choos-
166 ing the larger stimulus on the basis of r_1 and r_2 . One can then consider alternative *encoding rules*
167 for mapping objective stimulus magnitudes to feasible internal representations. The answer to this
168 *efficient coding* problem generally depends on the prior distribution $f(v)$ from which the different
169 stimulus magnitudes v_i are drawn. The resources required for more precise internal representa-
170 tions of individual stimuli may be economized with respect to either or both of two distinct cognitive
171 costs. The first goal of this work is to distinguish between these two types of efficiency concerns.

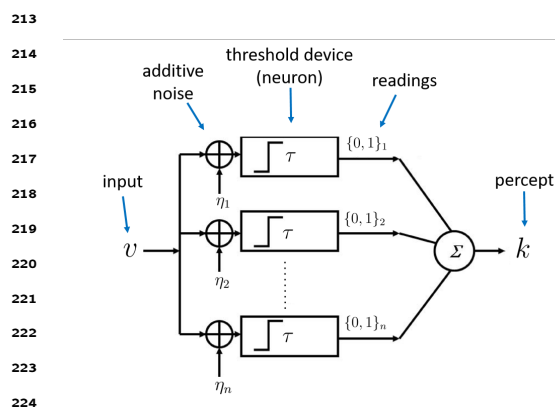
172 One question that we can ask is whether the observed behavioral responses are consistent with
173 the hypothesis that the conditional probabilities $p(r|v)$ are well-adapted to the particular frequency
174 distribution of stimuli used in the experiment, suggesting an efficient allocation of the limited en-
175 coding neural resources. The assumption of full adaptation is typically adopted in efficient coding
176 formulations of early sensory systems (*Laughlin, 1981; Wei and Stocker, 2017*), and also more re-
177 cently in applications of efficient coding theories in value-based decisions (*Louie and Glimcher,*
178 *2012; Polanía et al., 2019; Rustichini et al., 2017*).

179 There is also a second cost in which it may be important to economize on cognitive resources.
180 An efficient coding scheme in the sense described above economizes on the resources used to
181 represent each individual new stimulus that is encountered; however, the encoding and decoding
182 rules are assumed to be precisely optimized for the specific distribution $f(v)$ of stimuli that char-
183 acterizes the experimental situation. In practice, it will be necessary for a decision maker to learn
184 about this distribution in order to encode and decode individual stimuli in an efficient way, on the
185 basis of experience with a given context. In this case, the relevant design problem should not be
186 conceived as choosing conditional probabilities $p(r|v)$ once and for all, with knowledge of the prior
187 distribution $f(v)$ from which v will be drawn. Instead, it should be to choose a rule that specifies
188 how the probabilities $p(r|v)$ in the case of an individual stimulus should adapt to the distribution
189 of stimuli that have been encountered in a given context. It then becomes possible to consider
190 how well a given *learning* rule economizes on the degree of information about the distribution
191 of magnitudes associated with one's current context that is required for a given level of average
192 performance across contexts. This issue is important not only to reduce the cognitive resources
193 required to implement the rule in a given context (by not having to store or access so detailed a

194 description of the prior distribution), but in order to allow faster adaptation to a new context when
 195 the statistics of the environment can change unpredictably (Młynarski and Hermundstad, 2019).

196 Coding architecture

197 We now make the contrast between these two types of efficiency more concrete by considering
 198 a specific architecture for internal representations of sensory magnitudes. We suppose that the
 199 representation r_i of a given stimulus will consist of the output of a finite collection of n processing
 200 units, each of which has only two possible output states ("high" or "low" readings), as in the case of a
 201 simple perceptron. The probability that each of the units will be in one output state or the other can
 202 depend on the stimulus v_i that is presented. We further restrict the complexity of feasible encoding
 203 rules by supposing that the probability of a given unit being in the "high" state must be given
 204 by some function $\theta(v_i)$ that is the same for each of the individual units, rather than allowing the
 205 different units to coordinate in jointly representing the situation in some more complex way. We
 206 argue that the existence of multiple units operating in parallel effectively allows multiple repetitions
 207 of the same "experiment", but does not increase the complexity of the kind of test that can be
 208 performed. Note that we do not assume any unavoidable degree of stochasticity in the functioning
 209 of the individual units; it turns out that in our theory, it will be efficient for the units to be stochastic,
 210 but we do not assume that precise, deterministic functioning would be infeasible. Our resource
 211 limits are instead on the number of available units, the degree of differentiation of their output
 212 states, and the degree to which it is possible to differentiate the roles of distinct units.



225 **Figure 1.** Architecture of the sampling mechanism.
 226 Each processing unit receives noisy versions of the
 227 input v , where the noisy signals are i.i.d. additive
 228 random signals independent of v . The output of the
 229 neuron for each sample is "high" (one) reading if
 230 $v - \eta > \tau$ and zero otherwise. The noisy percept of
 the input is simply the sum of the outputs of each
 sample given by k .

231
 232 The efficient coding problem for a given en-
 233 vironment, specified by a particular prior distribution $f(v)$, will be to choose the encoding rule $\theta(v)$
 234 so as to allow an overall distribution of responses across trials that will be as accurate as possi-
 235 ble (according to criteria that we will elaborate further below). We can further suppose that each
 236 of the individual processing units is a threshold unit, that produces a "high" reading if and only
 237 if the value $v_i - \eta_i$ exceeds some threshold τ , where η_i is a random term drawn independently
 238 on each trial from some distribution f_η (Figure 1). The encoding function $\theta(v)$ can then be imple-
 239 mented by choice of an appropriate distribution f_η . This implementation requires that $\theta(v)$ be a
 non-decreasing function, as we shall assume.

Given such a mechanism, the internal repre-
 sentation r_i of the magnitude of an individ-
 ual stimulus v_i will be given by the collection of
 output states of the n processing units. A spec-
 ification of the function $\theta(v)$ then implies con-
 ditional probabilities for each of the 2^n possi-
 ble representations. Given our assumption of
 a symmetrical and parallel process, the num-
 ber k_i of units in the "high" state will be a suf-
 ficient statistic, containing all of the informa-
 tion about the true magnitude v_i that can be
 extracted from the internal representation. An
 optimal decoding rule will therefore be a func-
 tion only of k_i , and we can equivalently treat
 k_i (an integer between 0 and n) as the internal
 representation of the quantity v_i . The condi-
 tional probabilities of different internal repre-
 sentations are then

$$p(k_i | v_i) = \binom{n}{k} \theta(v_i)^{k_i} (1 - \theta(v_i))^{n-k_i}. \quad (1)$$

240 Limited cognitive resources

241 One measure of the cognitive resources required by such a system is the number n of processing
242 units that must produce an output each time an individual stimulus v_i is evaluated. We can consider
243 the optimal choice of f_n in order to maximize, for instance, average accuracy of responses in a given
244 environment $f(v)$, in the case of any bound n on the number of units that can be used to represent
245 each stimulus. But we can also consider the amount of information about the distribution $f(v)$
246 that must be used in order to decide how to encode a given stimulus v_i . If the system is to be able
247 to adapt to changing environments, it must determine the value of θ (the probability of a "high"
248 reading) as a function of both the current v_i and information about the distribution f , in a way
249 that must now be understood to apply across different potential contexts. This raises the issue of
250 how precisely the distribution f associated with the current context is represented for purposes
251 of such a calculation. A more precise representation of the prior (allowing greater sensitivity to
252 fine differences in priors) will presumably entail a greater resource cost or very long adaptation
253 periods.

254 We can quantify the precision with which the prior f is represented by supposing that it is
255 represented by a finite sample of m independent draws, $\tilde{v}_1, \dots, \tilde{v}_m$, from the prior (or more precisely,
256 from the set of previously experienced values, an empirical distribution that should after sufficient
257 experience provide a good approximation to the true distribution). We further assume that an
258 independent sample of m previously experienced values is used by each of the processing units
259 (**Figure 1**). Each of the n individual processing units is then in the "high" state with probability
260 $\theta(v_i; \tilde{v}_1, \dots, \tilde{v}_m)$. The complete internal representation of the stimulus v_i is then the collection of n
261 independent realizations of this binary-valued random variable. We may suppose that the resource
262 cost of an internal representation of this kind is an increasing function of both n and m .

263 This allows us to consider an efficient coding meta-problem, in which for any given values (n, m) ,
264 the function $\theta(v_i; \tilde{v}_1, \dots, \tilde{v}_m)$ is chosen so as to maximize some measure of average perceptual accu-
265 racy, where the average is now taken not only over the entire distribution of possible v_i occurring
266 under a given prior $f(v)$, but over some range of different possible priors for which the adaptive
267 coding scheme is to be optimized. We wish to consider how each of the two types of resource
268 constraint (a finite bound on n as opposed to a finite bound on m) affects the nature of the pre-
269 dicted imprecision in internal representations, under the assumption of a coding scheme that is
270 efficient in this generalized sense, and then ask whether we can tell in practice how tight each of
271 the resource constraints appears to be.

272 Efficient sampling for a known prior distribution

273 We first consider efficient coding in the case that there is no relevant constraint on the size of m ,
274 while n instead is bounded. In this case, we can assume that each time an individual stimulus v_i
275 must be encoded, a large enough sample of prior values is used to allow accurate recognition of
276 the distribution $f(v)$, and the problem reduces to a choice of a function $\theta(v)$ that is optimal for each
277 possible prior $f(v)$.

278 Maximizing mutual information

279 The nature of the resource-constrained problem to be optimized depends on the performance
280 measure that we use to determine the usefulness of a given encoding scheme. A common as-
281 sumption in the literature on efficient coding has been that the encoding scheme maximizes the
282 mutual information between the true stimulus magnitude and its internal representation (**Ganguli**
283 **and Simoncelli, 2014; Polanía et al., 2019; Wei and Stocker, 2015**). We start by characterizing the
284 optimal $\theta(v)$ for a given prior distribution $f(v)$, according to this criterion. It can be shown that for
285 large n , the mutual information between θ and k (hence the mutual information between v and k)
286 is maximized if the prior distribution \hat{f} over θ is Jeffreys' prior (**Clarke and Barron, 1994**)

$$\hat{f}(\theta) = \frac{1}{\pi\sqrt{\theta(1-\theta)}}, \quad (2)$$

287 also known as the arcsine distribution. Hence, the mapping $\theta(v)$ induces a prior distribution \hat{f} over
288 θ given by the arcsine distribution (**Figure 2a**, right panel). Based on this result, it can be shown
289 that the optimal encoding rule $\theta(v)$ that guarantees maximization of mutual information between
290 the random variable v and the noisy encoded percept k is given by (see **Appendix 1**)

$$\theta(v) = \left[\sin \left(\frac{\pi}{2} F(v) \right) \right]^2, \quad (3)$$

291 where $F(v)$ is the CDF of the prior distribution $f(v)$.

292 Accuracy maximization for a known prior distribution

293 So far, we have derived the optimal encoding rule to maximize mutual information, however, one
294 may ask what the implications are of such a theory for discrimination performance. This is impor-
295 tant to investigate given that achieving channel capacity does not necessarily imply that the goals
296 of the organism are also optimized (**Park and Pillow, 2017**). Independent of information maximiza-
297 tion assumptions, here we start from scratch and investigate what are the necessary conditions
298 for minimizing discrimination errors given the resource-constrained problem considered here. We
299 solve this problem for the case of two alternative forced choice tasks, where the average probability
300 of error is given by (see **Appendix 2**)

$$E[\text{error}] = \iint P_{\text{error}}[\theta(v_1), \theta(v_2)] \hat{f}(\theta_1) \hat{f}(\theta_2) d\theta_1 d\theta_2, \quad (4)$$

301 where $P_{\text{error}}[\cdot]$ represents the probability of erroneously choosing the alternative with the lowest
302 value v given a noisy percept k (assuming that the goal of the organism in any given trial is to
303 choose the alternative with the highest value). Here, we want to find the density function $\hat{f}(\theta)$ that
304 guarantees the smallest average error (Eq. 4). The solution to this problem is (**Appendix 2**)

$$\hat{f}(\theta) = \frac{1}{\pi \sqrt{\theta(1-\theta)}}, \quad (5)$$

305 which is exactly the same prior density function over θ that maximizes mutual information (Eq. 2).
306 Crucially, please note that we have obtained this expression based on minimizing the frequency of
307 erroneous choices and not the maximization of mutual information as a goal in itself. This provides
308 a further (and normative) justification for why maximizing mutual information under this coding
309 scheme is beneficial when the goal of the agent is to minimize discrimination errors (i.e., maximize
310 accuracy).

311 Optimal noise for a known prior distribution

312 Based on the coding architecture presented in **Figure 1**, the optimal encoding function $\theta(v)$ can then
313 be implemented by choice of an appropriate distribution f_η . It can be shown that discrimination
314 performance can be optimized by finding the optimal noise distribution f_η (**Appendix 3**) (**McDonnell**
315 **et al., 2007**)

$$f_\eta(v) = \frac{\pi}{2} \sin[\pi(1 - F(\tau - v))] f(\tau - v). \quad (6)$$

316 Remarkably, this result is independent of the number of samples n available to encode the input
317 variable, and generalizes to any prior distribution f (recall that F is defined as its cumulative density
318 function).

319 This result reveals three important aspects of neural function and decision behavior: First, it
320 makes explicit why a system that evolved to code information using a coding scheme of the kind
321 assumed in our framework must be necessarily noisy. That is, we do not attribute the randomness
322 of peoples' responses to a particular set of stimuli or decision problem to unavoidable randomness
323 of the hardware used to process the information. Instead, the relevant constraints are assumed
324 to be the limited set of output states for each neuron, the limited number of neurons, and the
325 requirement that the neurons operate in parallel (so that each one's output state must be statis-
326 tically independent of the others, conditional on the input stimulus). Given these constraints, we

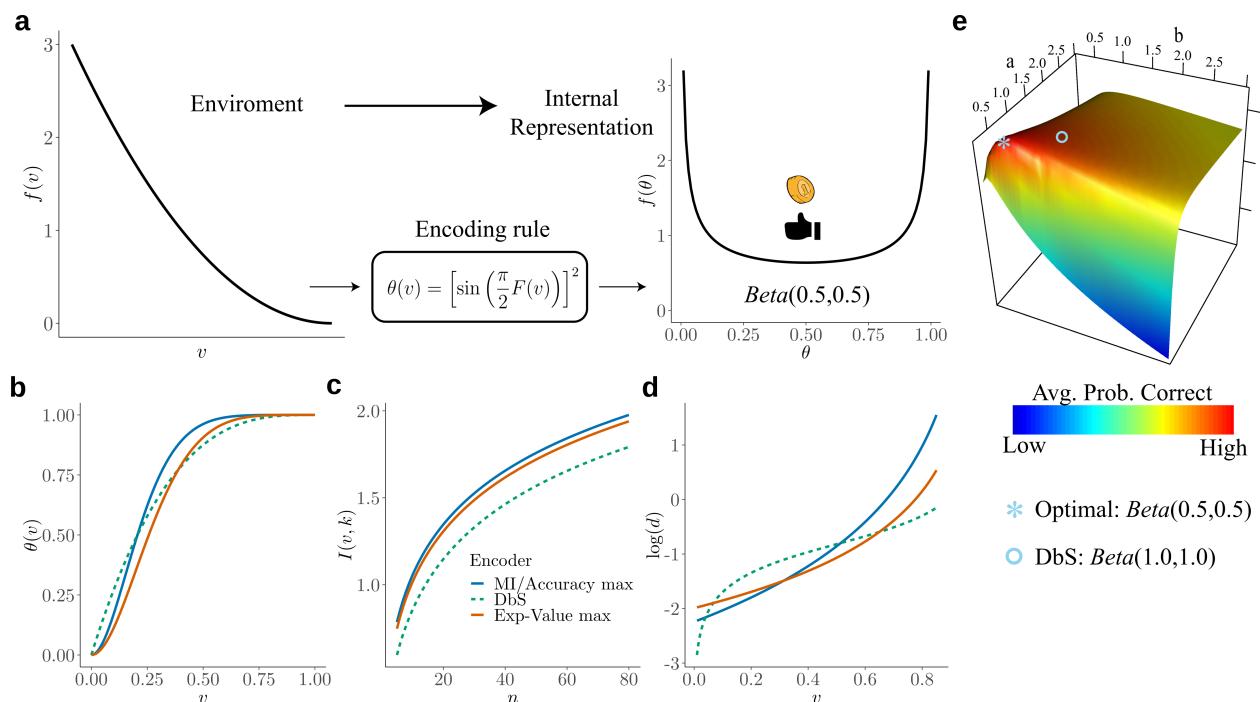


Figure 2. a) Schematic representation of our theory. Left: example prior distribution $f(v)$ of values v encountered in the environment. Right: Prior distribution in the encoder space (Eq. 2) due to optimal encoding (Eq. 3). This optimal mapping determines the probability θ of generating a "high" or "low" reading. The ex-ante distribution over θ that guarantees maximization of mutual information is given by the arcsine distribution (Eq. 2). **b)** Encoding rules $\theta(v)$ for different decision strategies under binary sampling coding: accuracy maximization (blue), reward maximization (red), DbS (green dashed). **c)** Mutual information $I(v, k)$ for the different encoding rules as a function of the number of samples n . As expected $I(v, k)$ increases with n , however the rule that results in the highest loss of information is DbS. **d)** Discriminability thresholds d (log-scaled for better visualization) for the different encoding rules as a function of the input values v for the prior $f(v)$ given in panel a. **e)** Graphical representation of the perceptual accuracy optimization landscape. We plot the average probability of correct responses for the large n limit using as benchmark a Beta distribution with parameters a and b. The blue star shows the average error probability assuming that $f(\theta)$ is the arcsine distribution (Eq. 2), which is the optimal solution when the prior distribution f is known. The blue open circle shows the average error probability based on the encoding rule assumed in DbS, which is located near the optimal solution. Please note that when formally solving this optimization problem, we did not assume a priori that the solution is related to the beta distribution. We use the beta distribution in this figure just as a benchmark for visualization. Detailed comparison of performance for finite n samples is presented in **Appendix 7**.

327 show that it is efficient for the operation of the neurons to be random. Second, it shows how the
 328 nervous system may take advantage of these noisy properties by reshaping its noise structure to
 329 optimize decision behavior. Third, it shows that the noise structure can remain unchanged irre-
 330 spective of the amount of resources available to guide behavior (i.e., the noise distribution f_{η} does
 331 not depend on n , Eq. 6). Please note however, that this minimalistic implementation does not
 332 directly imply that the samples in our algorithmic formulation are necessarily drawn in this way.
 333 We believe that this implementation provides a simple demonstration of the consequences of lim-
 334 ited resources in systems that encode information based on discrete stochastic events (*Sharpee,*
 335 **2017**). Interestingly, it has been shown that this minimalistic formulation can be extended to more
 336 realistic population coding specifications (*Nikitin et al., 2009*).

337 Efficient coding and the relation between environmental priors and discrimination
 338 The results presented above imply that this encoding framework imposes limitations on the abil-
 339 ity of capacity-limited systems to discriminate between different values of the encoded variables.
 340 Moreover, we have shown that error minimization in discrimination tasks implies a particular
 341 shape of the prior distribution of the encoder (Eq. 5) that is exactly the prior density that maxi-
 342 mizes mutual information between the input v and the encoded noisy readings k (Eq. 5, **Figure 2a**
 343 right panel). Does this imply a relation between prior and discriminability over the space of the

344 encoded variable? Intuitively, following the efficient coding hypothesis, the relation should be that
345 lower discrimination thresholds should occur for ranges of stimuli that occur more frequently in
346 the environment or context.

347 Recently, it was shown that using an efficiency principle for encoding sensory variables (e.g.,
348 with a heterogeneous population of noisy neurons (*Ganguli and Simoncelli, 2016*)) it is possible to
349 obtain an explicit relationship between the statistical properties of the environment and perceptual
350 discriminability (*Ganguli and Simoncelli, 2016*). The theoretical relation states that discriminability
351 thresholds d should be inversely proportional to the density of the prior distribution $f(v)$. Here,
352 we investigated whether this particular relation also emerges in the efficient coding scheme that
353 we propose in this study.

Remarkably, we obtain the following relation between discriminability thresholds, prior distribution of input variables, and the number of limited samples n (**Appendix 4**):

$$d = \frac{1}{\sqrt{n\pi}f(v)} \propto \frac{1}{f(v)} \quad (7)$$

354 Interestingly, this relationship between prior distribution and discriminability thresholds holds empirically
355 across several sensory modalities (**Appendix 4**), thus once again demonstrating that the
356 efficient coding framework that we propose here seems to incorporate the right kind of constraints
357 to explain observed perceptual phenomena as consequences of optimal allocation of finite capacity
358 for internal representations.

359 **Maximizing the expected size of the selected option (fitness maximization)**

360 Until now, we have studied the case when the goal of the organism is to minimize the number of
361 mistakes in discrimination tasks. However, it is important to consider the case when the goal of
362 the organism is to maximize fitness or expected reward (*Pirrone et al., 2014*). For example, when
363 spending the day foraging fruit, one must make successive decisions about which tree has more
364 fruits. Fitness depends on the number of fruit collected which is not a linear function of the number
365 of accurate decisions, as each choice yields a different amount of fruit.

366 Therefore, in the case of reward maximization, we are interested in minimizing reward loss
367 which is given by the following expression

$$E[v(\text{chosen})] = \int \int f(v_1, v_2) [P_1(\theta(v_1), \theta(v_2))v_1 + P_2(\theta(v_1), \theta(v_2))v_2] dv_1 dv_2, \quad (8)$$

368 where $P_i(\theta(v_1), \theta(v_2))$ is the probability of choosing option i when the input values are v_1 and v_2 .
369 Thus, the goal is to find the encoding rule $\theta(v)$ which guarantees that the amount of reward loss is
370 as small as possible given our proposed coding framework.

371 Here we show that the optimal encoding rule $\theta(v)$ that guarantees maximization of expected
372 value is given by

$$\theta(v) = \sin \left[\frac{\pi}{2} \cdot c \int_{-\infty}^v f(\tilde{v})^{2/3} d\tilde{v} \right]^2, \quad (9)$$

373 where c is a normalizing constant which guarantees that the expression within the integral is a probability
374 density function (**Appendix 5**). The first observation based on this result is that the encoding
375 rule for maximizing fitness is different from the encoding rule that maximizes accuracy (compare
376 Eqs. 3 and 9), which leads to a slight loss of information transmission (**Figure 2c**). Additionally, one
377 can also obtain discriminability threshold predictions for this new encoding rule. Assuming a right-
378 skewed prior distribution, which is often the case for various natural priors in the environment
379 (e.g., like the one shown in **Figure 2a**), we find that discriminability for small input values is lower
380 for reward maximization compared to perceptual maximization, however this pattern inverts for
381 higher values (**Figure 2d**). In other words, when we intend to maximize reward (given the shape

382 of our assumed prior, **Figure 2a**), the agent should allocate more resources to higher values (com-
383 pared to the perceptual case), however without completely giving up sensitivity for lower values,
384 as these values are still encountered more often.

385 **Efficient sampling with costs on acquiring prior knowledge**

386 In the previous section, we obtained analytical solutions that approximately characterize the opti-
387 mal $\theta(v)$ in the limit as n is made sufficiently large. Note however that we are always assuming
388 that n is finite, and that this constrains the accuracy of the decision maker's judgments, while m is
389 instead unbounded and hence no constraint.

390 The nature of the optimal function $\theta(v; \tilde{v}_1, \dots, \tilde{v}_m)$ is different, however, when m is small. We
391 argue that this scenario is particularly relevant when full knowledge of the prior is not warranted
392 given the costs vs benefits of learning, for instance, when the system expects contextual changes
393 to occur often. In this case, as we will formally elaborate below, it ceases to be efficient for θ to vary
394 only gradually as a function of v_i , rather than moving abruptly from values near zero to values near
395 one (**Appendix 6**). In the large- m limiting case, the distributions of sample values $(\tilde{v}_1, \dots, \tilde{v}_m)$ used
396 by the different processing units will be nearly the same for each unit (approximating the current
397 true distribution $f(v)$). Then if θ were to take only the values zero and one for different values of
398 its arguments, the n units would simply produce n copies of the same output (either zero or one)
399 for any given stimulus v_i and distribution $f(v)$. Hence only a very coarse degree of differentiation
400 among different stimulus magnitudes would be possible. Having θ vary more gradually over the
401 range of values of v_i in the support of $f(v)$ instead makes the representation more informative. But
402 when m is small (e.g., because of costs vs benefits of accurately representing the prior f), this kind
403 of arbitrary randomization in the output of individual processing units is no longer essential. There
404 will already be considerable variation in the outputs of the different units, even when the output of
405 each unit is a deterministic function of $(v_i; \tilde{v}_1, \dots, \tilde{v}_m)$, owing to the variability in the sample of prior
406 observations that is used to assess the nature of the current environment. As we will show below,
407 this variability will already serve to allow the collective output of the several units to differentiate
408 between many gradations in the magnitude of v_i , rather than only being able to classify it as "small"
409 or "large" (because either all units are in the "low" or "high" states).

410 **Robust optimality of Decision by Sampling**

411 Because of the way in which sampling variability in the values $(\tilde{v}_1, \dots, \tilde{v}_m)$ used to adapt each unit's
412 encoding rule to the current context can substitute for the arbitrary randomization represented
413 by the noise term η_i (see **Figure 1**), a sharp reduction in the value of m need not involve a great loss
414 in performance relative to what would be possible (for the same limit on n) if m were allowed to be
415 unboundedly large (**Appendix 7**). As an example, consider the case in which $m = 1$, so that each unit
416 j 's output state must depend only the value of the current stimulus v_i and one randomly selected
417 draw \tilde{v}_j from the prior distribution $f(v)$. A possible decision rule that is radically economical in this
418 way is one that specifies that the unit will be in the "high" state if and only if $v_i > \tilde{v}_j$. In this case,
419 the internal representation of a stimulus v_i will be given by the number k_i out of n independent
420 draws from the contextual distribution $f(v)$ with the property that the contextual draw is smaller
421 than v_i , as in the model of *decision by sampling* (DbS) (**Stewart et al., 2006**). However, it remains to
422 be determined to what degree it might be beneficial for a system to adopt such coding strategy.

423 In any given environment (characterized by a particular contextual distribution $f(v)$), DbS will be
424 equivalent to an encoding process with an architecture of the kind shown in **Figure 1**, but in which
425 the distribution $f_\eta = f(v)$ (compare to the optimal noise distribution f_η for the full prior adaptation
426 case in Eq. 6). This makes $\theta(v)$ vary endogenously depending on the contextual distribution $f(v)$.
427 And indeed, the way that $\theta(v)$ varies with the contextual distribution under DbS is fairly similar to
428 the way in which it would be optimal for it to vary in the absence of any cost of precisely learn-
429 ing and representing the contextual distribution. This result implies that $\theta(v)$ will be a monotonic
430 transformation of a function that increases more steeply over those regions of the stimulus space

431 where $f(v)$ is higher, regardless of the nature of the contextual distribution. We consider its perfor-
432 mance in a given environment, from the standpoint of each of the possible performance criteria
433 considered for the case of full prior adaptation (i.e., maximize accuracy or fitness), and show that it
434 differs from the optimal encoding rules under any of those criteria (**Figure 2b-d**). In particular, here
435 we show that using the encoding rule employed in DbS results in considerable loss of information
436 compared to the full-prior adaptation solutions (**Figure 2c**). An additional interesting observation
437 is that for the strategy employed in DbS, the agent appears to be more sensitive for extreme input
438 values, at least for a wide set of skewed distributions (e.g., for the prior distribution $f(v)$ in **Fig-**
439 **ure 2a**, the discriminability thresholds are lower at the extremes of the support of $f(v)$). In other
440 words, agents appear to be more sensitive to salience in the DbS rule. Despite these differences,
441 here it is important to emphasize that in general for all optimization objectives, the encoding rules
442 will be steeper for regions of the prior with higher density. However, mild changes in the steepness
443 of the curves will be represented in significant discriminability differences between the different
444 encoding rules across the support of the prior distribution (**Figure 2d**).

445 While the predictions of DbS are not exactly the same as those of efficient coding in the case
446 of unbounded m , under any of the different objectives that we consider, our numerical results
447 show that it can achieve performance nearly as high as that of the theoretically optimal encoding
448 rule; hence radically reducing the value of m does not have a large cost in terms of the accuracy of
449 the decisions that can be made using such an internal representation (**Appendix 7** and **Figure 2e**).
450 Under the assumption that reducing either m or n would serve to economize on scarce cognitive
451 resources, we formally prove that it might well be most efficient to use an algorithm with a very
452 low value of m (even $m = 1$, as assumed by DbS), while allowing n to be much larger (**Appendix 6**,
453 **Appendix 7**).

454 Crucially, here it is essential to emphasize that the above-mentioned results are derived for the
455 case of a particular finite number of processing units n (and a corresponding finite total number
456 of samples from the contextual distribution used to encode a given stimulus), and do not require
457 that n must be large (**Appendix 6**, **Appendix 7**).

458 **Testing theories of numerosity discrimination**

459 Our goal now is to compare back-to-back the resource-limited coding frameworks elaborated above
460 in a fundamental cognitive function for human behavior: numerosity perception. We designed
461 a set of experiments that allowed us to test whether human participants would adapt their nu-
462 merosity encoding system to maximize fitness or accuracy rates via full prior adaptation as usually
463 assumed in optimal models, or whether humans employ a "less optimal" but more efficient strat-
464 egy such as DbS, or the more established logarithmic encoding model.

465 In Experiment 1, healthy volunteers ($n=7$) took part in a two-alternative forced choice numeros-
466 ity task, where each participant completed $\sim 2,400$ trials across four consecutive days (methods).
467 On each trial, they were simultaneously presented with two clouds of dots and asked which one
468 contained more dots, and were given feedback on their reward and opportunity losses on each
469 trial (**Figure 3a**). Participants were either rewarded for their accuracy (perceptual condition, where
470 maximizing the amount of correct responses is the optimal strategy) or the number of dots they
471 selected (value condition, where maximizing reward is the optimal strategy). Each condition was
472 tested for two consecutive days with the starting condition randomized across participants. Cru-
473 cially, we imposed a prior distribution $f(v)$ with a right-skewed quadratic shape (**Figure 3b**), whose
474 parametrization allowed tractable analytical solutions of the encoding rules $\theta_A(v)$, $\theta_R(v)$ and $\theta_D(v)$,
475 that correspond to the encoding rules for Accuracy maximization, Reward maximization, and DbS,
476 respectively (**Figure 3e** and methods). Qualitative predictions of behavioral performance indicate
477 that the accuracy maximization model is the most accurate for trials with lower numerosities (the
478 most frequent ones), while the reward-maximization model outperforms the others for trials with
479 larger numerosities (trials where the difference in the number of dots in the clouds, and thus the
480 potential reward, is the largest, **Figure 2d** and **Figure 3f**). In contrast, the DbS strategy presents

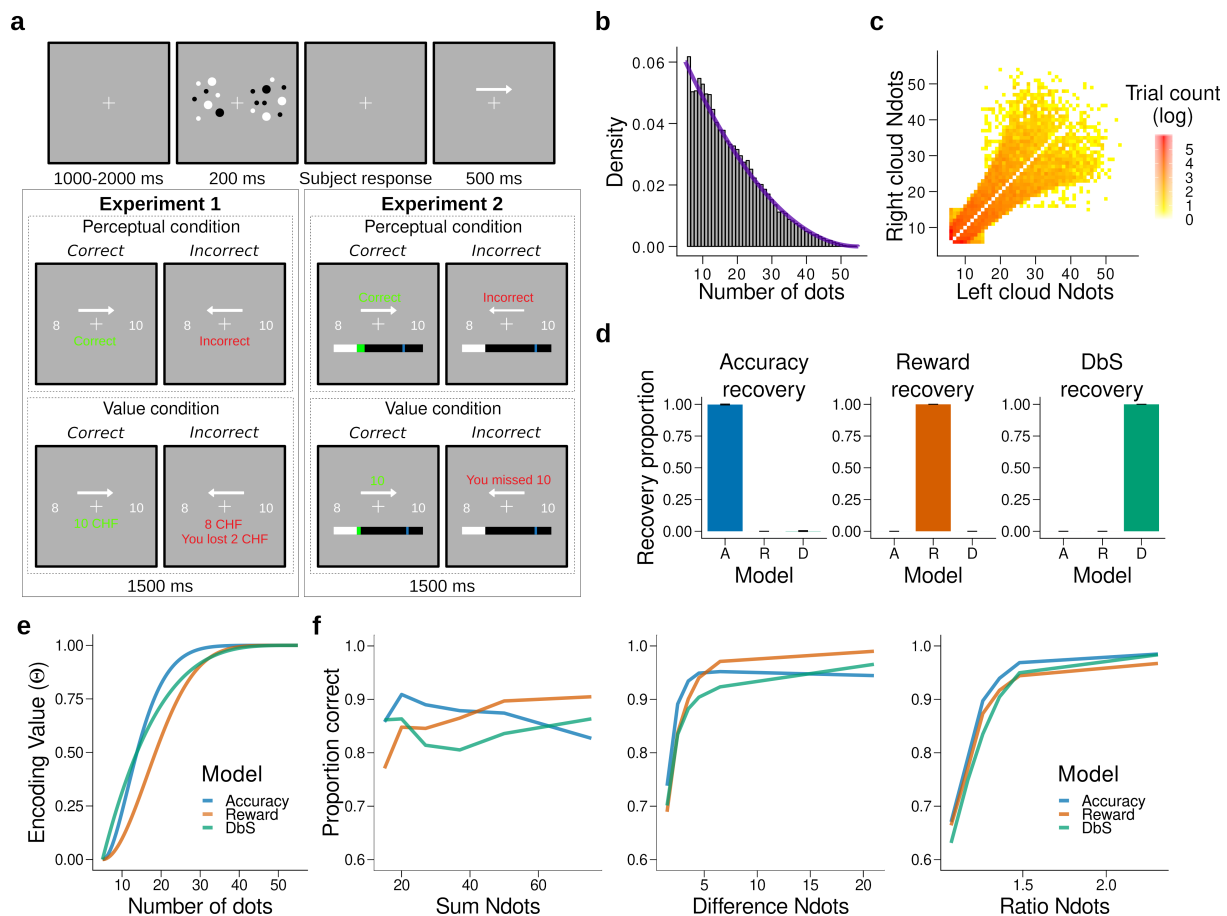


Figure 3. Experimental design, model simulations and recovery. **a**) Schematic task design of Experiments 1 and 2. After a fixation period (1-2s) participants were presented two clouds of dots (200ms) and had to indicate which cloud contained the most dots. Participants were rewarded for being accurate (Perceptual condition) or for the number of dots they selected (Value condition) and were given feedback. In Experiment 2 participants collected on correctly answered trials a number of points equal to a fixed amount (Perception condition) or a number equal to the dots in the cloud they selected (Value condition) and had to reach a threshold of points on each run. **b**) Empirical (grey bars) and theoretical (purple line) distribution of the number of dots in the clouds of dots presented across Experiments 1 and 2. **c**) Distribution of the numerosity pairs selected per trial. **d**) Synthetic data preserving the trial set statistics and number of trials per participant used in Experiment 1 was generated for each encoding rule (Accuracy (left), Reward (middle), and DbS (right)) and then the latent-mixture model was fit to each generated dataset. The figures show that it is theoretically possible to recover each generated encoding rule. **e**) Encoding function $\theta(v)$ for the different sampling strategies as a function of the input values v (i.e., the number of dots). **f**) Qualitative predictions of the three models (blue: Accuracy, red: Reward, green: Decision by Sampling) on trials from experiment 1 with $n = 25$. Performance of each model as a function of the sum of the number of dots in both clouds (left), the absolute difference between the number of dots in both clouds (middle) and the ratio of the number of dots in the most numerous cloud over the less numerous cloud (right).

Figure 3-Figure supplement 1. Model recovery for α fixed.

Figure 3-Figure supplement 2. Discriminability differences between the different encoding rules.

Figure 3-Figure supplement 3. Model recovery with both α and n as free parameters.

481 markedly different performance predictions, in line with the discriminability predictions of our formal analyses (Figure 2c,d).

482
483 In our modelling specification, the choice structure is identical for the three different sampling
484 models, differing only in the encoding rule $\theta(v)$ (methods). Therefore, answering the question of
485 which encoding rule is the most favored for each participant can be parsimoniously addressed
486 using a latent-mixture model, where each subject uses $\theta_A(v)$, $\theta_R(v)$ or $\theta_D(v)$ to guide their decisions
487 (methods). Before fitting this model to the empirical data, we confirmed the validity of our model
488 selection approach through a validation procedure using synthetic choice data (Figure 3d, Figure 3-
489 Figure Supplement 1, and methods).

490 After we confirmed that we can reliably differentiate between our competing encoding rules,
491 the latent-mixture model was initially fit to each condition (perceptual or value) using a hierarchi-
492 cal Bayesian approach (methods). Surprisingly, we found that participants did not follow the ac-
493 curacy or reward optimization strategy in the respective experimental condition, but favored the
494 DbS strategy (proportion that DbS was deemed best in the perceptual $p_{\text{DbSfavored}} = 0.86$ and value
495 $p_{\text{DbSfavored}} = 0.93$ conditions, **Figure 4**). Importantly, this population-level result also holds at the in-
496 dividual level: DbS was strongly favored in 6 out of 7 participants in the perceptual condition, and
497 7 out of 7 in the value condition (**Figure 4–Figure Supplement 1**). These results are not likely to be
498 affected by changes in performance over time, as performance was stable across the four consecu-
499 tive days (**Figure 4–Figure Supplement 2**). Additionally, we investigated whether biases induced by
500 choice history effects may have influenced our results (**Abrahamyan et al., 2016; Keung et al., 2019;**
501 **Talluri et al., 2018**). Therefore, we incorporated both choice- and correctness-dependence history
502 biases in our models and fitted the models once again (methods). We found similar results to the
503 history-free models ($p_{\text{DbSfavored}} = 0.87$ in accuracy and $p_{\text{DbSfavored}} = 0.93$ in value conditions, **Figure 4c**).
504 At the individual level, DbS was again strongly favored in 6 out of 7 participants in the perceptual
505 condition, and 7 out of 7 in the value condition (**Figure 4–Figure Supplement 1**).

506 In order to investigate further the robustness of this effect, we introduced a slight variation in
507 the behavioral paradigm. In this new experiment (Experiment 2), participants were given points on
508 each trial and had to reach a certain threshold in each run for it to be eligible for reward (**Figure 3a**
509 and methods). This class of behavioral task is thought to be in some cases more ecologically valid
510 than trial-independent choice paradigms (**Kolling et al., 2014**). In this new experiment, either a
511 fixed amount of points for a correct trial was given (perceptual condition) or an amount equal to
512 the number of dots in the chosen cloud if the response was correct (value condition). We recruited
513 a new set of participants ($n=6$), who were tested on these two conditions, each for two consecutive
514 days with the starting condition randomized across participants (each participant completed 2,560
515 trials). The quantitative results revealed once again that participants did not change their encoding
516 strategy depending on the goals of the task, with DbS being strongly favored for both perceptual
517 and value conditions ($p_{\text{DbSfavored}} = 0.999$ and $p_{\text{DbSfavored}} = 0.91$, respectively; **Figure 4a**), and these
518 results were confirmed at the individual level where DbS was strongly favored in 6 out of 6 partici-
519 pants in both the perceptual and value conditions (**Figure 4–Figure Supplement 1**). Once again,
520 we found that inclusion of choice history biases in this experiment did not significantly affect our
521 results both at the population and individual levels. Population probability that DbS was deemed
522 best in the perceptual ($p_{\text{DbSfavored}} = 0.999$) and value ($p_{\text{DbSfavored}} = 0.90$) conditions (**Figure 4–Figure**
523 **Supplement 1**), and at the individual level DbS was strongly favored in 6 out of 6 participants in
524 the perceptual condition and 5 of 6 in the value condition (**Figure 4–Figure Supplement 1**). Thus,
525 experiments 1 and 2 strongly suggest that our results are not driven by specific instructions or
526 characteristics of the behavioral task.

527 As a further robustness check, for each participant we grouped the data in different ways across
528 experiments (Experiments 1 and 2) and experimental conditions (perceptual or value) and investi-
529 gated which sampling model was favored. We found that irrespective of how the data was grouped,
530 DbS was the model that was clearly deemed best at the population (**Figure 4**) and individual level
531 (**Figure 4–Figure Supplement 3**). Additionally, we investigated whether these quantitative results
532 specifically depended on our choice of using a latent-mixture model. Therefore, we also fitted each
533 model independently and compared the quality of the model fits based on out-of-sample cross-
534 validation metrics (methods). Once again, we found that the DbS was favored independently of
535 experiment and conditions (**Figure 4**).

536 One possible reason why the two experimental conditions did not lead to differences could be
537 that, after doing one condition for two days, the participants did not adapt as easily to the new
538 incentive rule. However, note that as the participants did not know of the second condition before
539 carrying it out, they could not adopt a compromise between the two behavioral objectives. Never-
540 theless, we fitted the latent-mixture model only to the first condition that was carried out by each

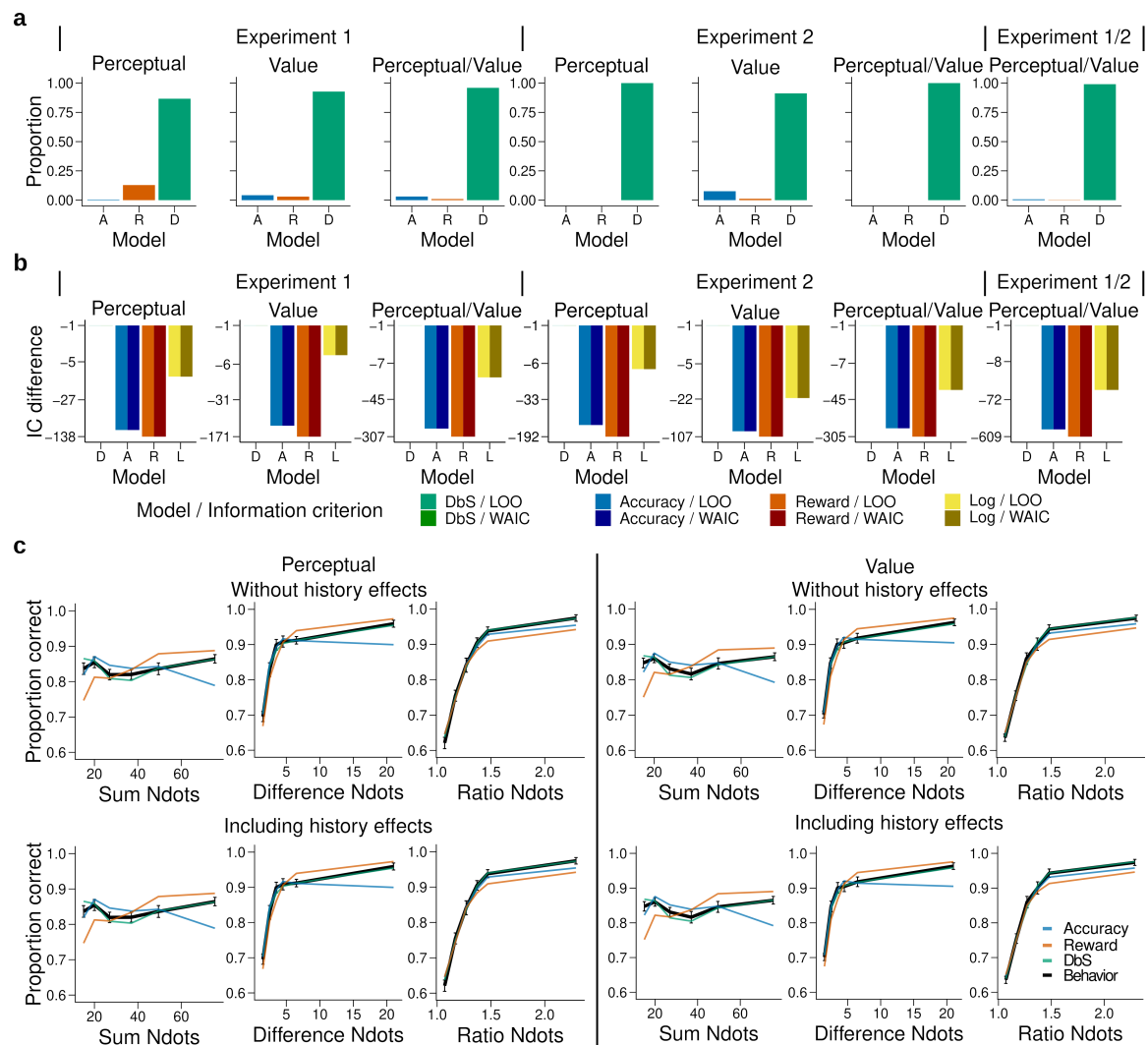


Figure 4. Behavioral results. **a**) Bars represent proportion of times an encoding rule (Accuracy (A, blue), Reward (R, red), DbS (D, green)) was selected by the Bayesian latent-mixture model based on the posterior estimates across participants. Each panel shows the data grouped for each and across experiments and experimental conditions (see titles on top of each panel). The results show that DbS was clearly the favored encoding rule. The latent vector π posterior estimates are presented in **Figure 4–Figure Supplement 8**. **b**) Difference in LOO and WAIC between the best model (DbS (D) in all cases) and the competing models: Accuracy (A), Reward (R) and Logarithmic (L) models. Each panel shows the data grouped for each and across experimental conditions and experiments (see titles on top of each panel). **c**) Behavioral data (black, error bars represent s.e.m. across participants) and model predictions based on fits to the empirical data. Data and model predictions are presented for both the perceptual (left panels) or value (right panels) conditions, and excluding (top panels) or including (bottom) choice history effects. Performance of data model predictions are presented as function of the sum of the number of dots in both clouds (left), the absolute difference between the number of dots in both clouds (middle) and the ratio of the number of dots in the most numerous cloud over the less numerous cloud (right). Results reveal a remarkable overlap of the behavioral data and predictions by DbS, thus confirming the quantitative results presented in panels a and b.

Figure 4–Figure supplement 1. Latent mixture model fits for each participant.

Figure 4–Figure supplement 2. Performance across time.

Figure 4–Figure supplement 3. Individual level fit of the latent mixture model combining data across experiments and experimental conditions.

Figure 4–Figure supplement 4. Model comparison based on leave-one-out cross-validation metrics.

Figure 4–Figure supplement 5. Reaction times are similar in the perceptual and value conditions.

Figure 4–Figure supplement 6. Behavior and model predictions as a function of sum and difference in dots.

Figure 4–Figure supplement 7. Model fit for the first experimental condition of each participant.

Figure 4–Figure supplement 8. Latent vector π posterior estimates.

541 participant. We found once again that DbS was the best model explaining the data, irrespective
542 of condition and experimental paradigm (**Figure 4–Figure Supplement 7**). Therefore, the fact that
543 DbS is favored in the results is not an artifact of carrying out two different conditions in the same
544 participants.

545 We also investigated whether the DbS model makes more accurate predictions than the widely
546 used logarithmic model of numerosity discrimination tasks (*Dehaene, 2003*). We found that DbS
547 still made better out of sample predictions than the log-model (**Figure 4b, Figure 5f,g**). Moreover,
548 these results continued to hold after taking into account possible choice history biases (**Figure 4–**
549 **Figure Supplement 4**). In addition to these quantitative results, qualitatively we also found that
550 behavior closely matched the predictions of the DbS model remarkably well (**Figure 4c**), based on
551 virtually only 1 free parameter, namely, the number of samples (resources) n . Together, these re-
552 sults provide compelling evidence that DbS is the most likely resource-constrained sampling strat-
553 egy used by participants in numerosity discrimination tasks.

554 Recent studies have also investigated behavior in tasks where perceptual and preferential deci-
555 sions have been investigated in paradigms with identical visual stimuli (*Dutilh and Rieskamp, 2016;*
556 *Polanía et al., 2014; Grueschow et al., 2015*). In these tasks, investigators have reported differences
557 in behavior, in particular in the reaction times of the responses, possibly reflecting differences in
558 behavioral strategies between perceptual and value-based decisions. Therefore, we investigated
559 whether this was the case also in our data. We found that reaction times did not differ between
560 experimental conditions for any of the different performance assessments considered here (**Fig-**
561 **ure 4–Figure Supplement 5**). This further supports the idea that subjects were in fact using the
562 same sampling mechanism irrespective of behavioral goals.

563 Here it is important to emphasize that all sampling models and the logarithmic model of nu-
564 merosity have the same degrees of freedom (performance is determined by n in the sampling
565 models and Weber’s fraction σ in the log model, methods). Therefore, qualitative and quantitative
566 differences favoring the DbS model cannot be explained by differences in model complexity. It
567 could also be argued that normal approximation of the binomial distributions in the sampling de-
568 cision models only holds for large enough n . However, we find evidence that the large- n optimal
569 solutions are also nearly optimal for low n values (**Appendix 7**). Estimates of n in our data are in
570 general $n \approx 21$ (**Table 1**) and we find that the large- n rule is nearly optimal already for $n = 15$ (**Ap-**
571 **pendix 7**). Therefore the asymptotic approximations should not greatly affect the conclusions of
572 our work.

573 Dynamics of adaptation

574 Up to now, fits and comparison across models have been done under the assumption that the
575 participants learned the prior distribution $f(v)$ imposed in our task. If participants are employing
576 DbS, it is important to understand the dynamical nature of adaptation in our task. Note that the
577 shape of the prior distribution is determined by the parameter α (**Figure 5b**, Eq. 10 in methods).
578 First, we made sure based on model recovery analyses that the DbS model could jointly and accu-
579 rately recover both the shape parameter α and the resource parameter n based on synthetic data
580 (**Figure 3–Figure Supplement 3**). Then we fitted this model to the empirical data and found that
581 the recovered value of the shape parameter α closely followed the value of the empirical prior with
582 a slight underestimation (**Figure 5a**). Next, we investigated the dynamics of prior adaptation. To
583 this end, we ran a new experiment (Experiment 3, $n=7$ new participants) where we set the shape
584 parameter of the prior to a lower value compared to Experiments 1-2 (**Figure 5b**, methods). We
585 investigated the change of α over time by allowing this parameter to change with trial experience
586 (Eq. 18, methods) and compared the evolution of α for Experiments 1 and 2 (empirical $\alpha = 2$) with
587 Experiment 3 (empirical $\alpha = 1$, **Figure 5b**). If participants show prior adaption in our numerosity dis-
588 crimination task, we hypothesized that the asymptotic value of α should be higher for Experiments
589 1-2 than for Experiment 3. First, we found that for Experiments 1-2, the value of α quickly reached
590 an asymptotic value close to the target value (**Figure 5c**). On the other hand, for Experiment 3

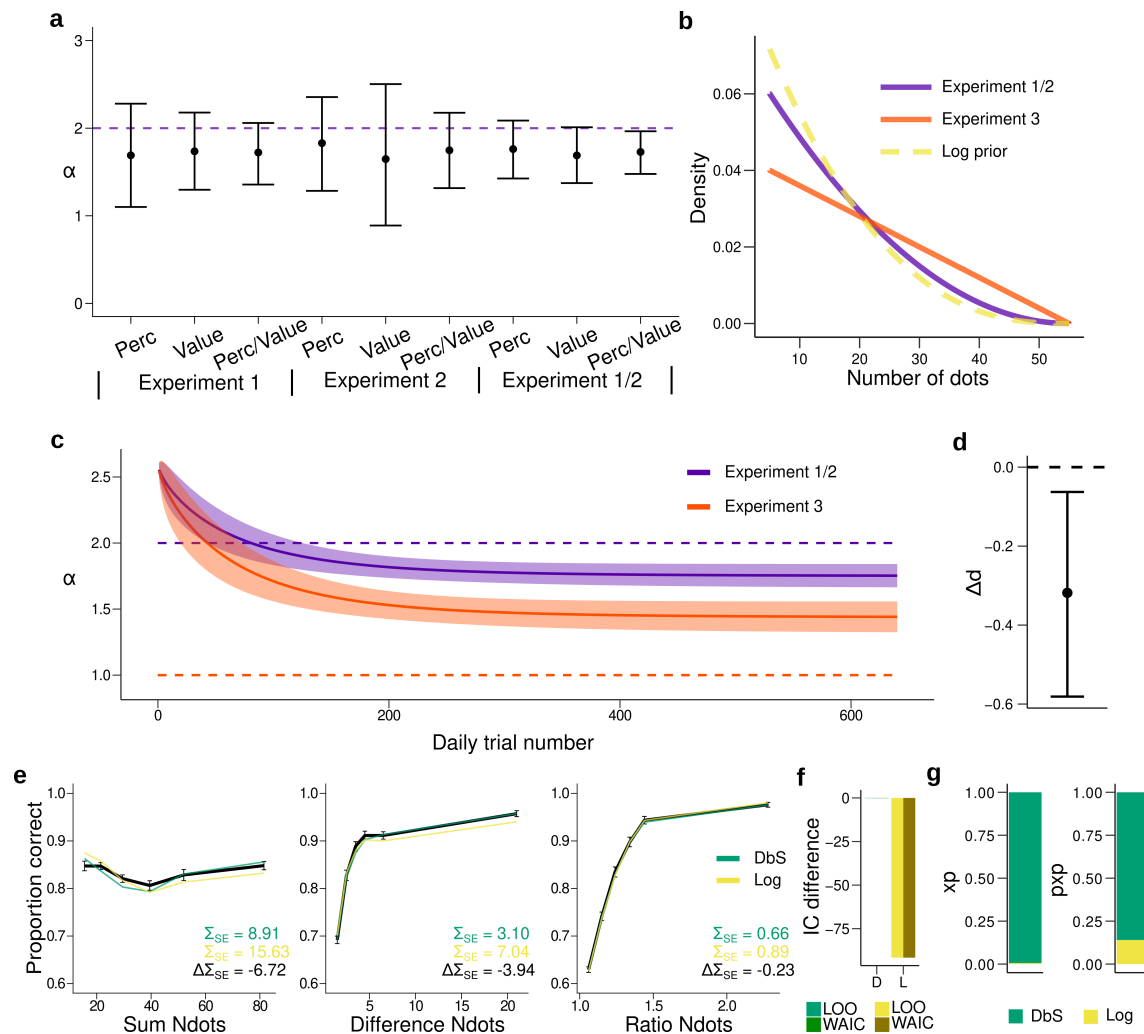


Figure 5. Prior adaptation analyses. **a)** Estimation of the shape parameter α for the DbS model by grouping the data for each and across experimental conditions and experiments. Error bars represent the 95% highest density interval of the posterior estimate of α at the population level. The dashed line shows the theoretical value of α . **b)** Theoretical prior distribution $f(v)$ in Experiments 1 and 2 ($\alpha = 2$, purple) and 3 ($\alpha = 1$, orange). The dashed line represents the value of α of our prior parametrization that approximates the DbS and log discriminability models. **c)** Posterior estimation of α_t (Eq. 18) as a function of the number of trials t in each daily session for Experiments 1 and 2 (purple) and Experiment 3 (orange). The results reveal that, as expected, α_t reaches a lower asymptotic value δ . Error bars represent ± 1 s.d. of 3,000 simulated α_t values drawn from the posterior estimates of the HBM (see methods). **d)** Difference in the δ parameter between Experiments 1-2 and Experiment 3 based on the posterior parameter estimates of the HBM. This analysis reveals a significant difference ($P_{\text{MCMC}} = 0.006$). Error bars represent the 95% highest density interval of the posterior differences in the HBM. **e)** Behavioral data (black) and model fit predictions of the DbS (green) and Log (yellow) models. Performance of each model as a function of the sum of the number of dots in both clouds (left), the absolute difference between the number of dots in both clouds (middle) and the ratio of the number of dots in the most numerous cloud over the less numerous cloud (right). Error bars represent s.e.m. **f)** Difference in LOO and WAIC between the best fitting DbS (D) and logarithmic encoding (Log) model. **g)** Population exceedance probabilities (xp, left) and protected exceedance probabilities (pxp, right) for DbS (green) vs Log (yellow) of a Bayesian model selection analysis (Stephan et al., 2009): $x_{\text{pDbS}} = 0.99$, $\text{px}_{\text{pDbS}} = 0.87$. These results provide a clear indication that the adaptive DbS explains the data better than the Log model.

Figure 5-Figure supplement 1. Performance across trial experience.

Figure 5-Figure supplement 2. Quantitative and dynamical analysis of adaptation over time

591 the value of α continued to decrease during the experimental session, but slowly approaching its
 592 target value. This seemingly slower adaptation to the shape of the prior in Experiment 3 might
 593 be explained by the following observation. The prior parametrized with $\alpha = 1$ in Experiment 3, is
 594 further away from an agent hypothesized to have a natural numerosity discrimination based on
 595 a log scale ($\alpha = 2.58$, Figure 5b and methods), which is closer in value to the shape of the prior in

596 Experiments 1 and 2 ($\alpha = 2$). This result is in line with previous DbS studies showing that adaptation
597 to "unnatural" priors in decision tasks is slower (*Stewart et al., 2015*) and may require many more
598 trials or training experience than it is possible in laboratory experiments. Irrespective of these con-
599 siderations, the key result to confirm our adaptation hypothesis is that the asymptotic value of α is
600 lower for Experiment 3 compared to Experiments 1 and 2 ($P_{\text{MCMC}} = 0.006$; *Figure 5c,d*). Additionally,
601 we found that this DbS model again provides more accurate qualitative and quantitative out of
602 sample predictions than the log model (*Figure 5e,f*).

603 We further investigated evidence for adaptation using an alternative quantitative approach.
604 First, we performed out of sample model comparisons based on the following models: (i) the
605 adaptive- α model, (ii) free- α model with alpha free but non-adapting over time, and (iii) fixed- α
606 model with $\alpha = 2$. The results of the out of sample predictions revealed that the best model was
607 the free- α model, followed closely by the the adaptive- α model ($\Delta\text{LOO} = 1.8$) and then by fixed- α
608 model ($\Delta\text{LOO} = 32.6$). However, we did not interpret the apparent small difference between the
609 adaptive- α and free- α model as evidence for lack of adaptation, given that the more complex adap-
610 tive model will be strongly penalized after adaptation is stable. That is, if adaptation is occurring,
611 then the adaptive- α only provides a better fit for the trials corresponding to the adaptation period.
612 After adaptation the adaptive- α should provide a similar fit than the free- α model, however with
613 a larger complexity that will be penalized by model comparison metrics. Therefore, to investigate
614 the presence of adaptation, we took a closer quantitative look at evolution of the fits across trial
615 experience. We computed the average trial-wise predicted Log-Likelihood (by sampling from the
616 hierarchical Bayesian model) and compared the differences of this metric between the competing
617 models and the adaptive model. We hypothesized that if adaptation is taking place, the adaptive- α
618 model would have an advantage relative to the free- α model at the beginning of the session, with
619 these differences vanishing towards the end. On the other hand, the fixed- α should roughly match
620 the adaptive- α model at the beginning and then become worse over time, but these differences
621 should stabilize after the end of the adaptation period. The results of these analyses support our
622 hypotheses (*Figure 5–Figure Supplement 2*), thus providing further evidence of adaptation, high-
623 lighting the fact that the DbS model can parsimoniously capture adaptation contextual changes in
624 a continuous and dynamical manner.

625 Discussion

626 The brain is a metabolically expensive inference machine (*Hawkes et al., 1998; Navarrete et al.,*
627 *2011*). Therefore it has been suggested that evolutionary pressure has driven it to make produc-
628 tive use of its limited resources by exploiting statistical regularities (*Attneave, 1954; Laughlin, 1981*).
629 Here, we incorporate this important – often ignored – aspect in models of behavior by introducing
630 a general framework of decision-making under the constraints that the system: (i) encodes infor-
631 mation based on binary codes, (ii) has limited number of samples available to encode information,
632 and (iii) considers the costs of both contextual adaptation.

633 Under the assumption that the organism has fully adapted to the statistics in a given context,
634 we show that the encoding rule that maximizes mutual information is the same rule that maximizes
635 decision accuracy in two-alternative decision tasks. However, note that there is nothing privileged
636 about maximizing mutual information, as it does not mean that the goals of the organism are
637 necessarily achieved (*Park and Pillow, 2017*). In fact, we show that if the goal of the organism is
638 instead to maximize the expected value of the chosen options, the system should not rely on max-
639 imizing information transmission to fulfill this goal and must give up a small fraction of precision
640 in information coding. Here, we derived analytical solution for each of these optimization objec-
641 tive criteria, emphasizing that these analytical solutions were derived for the large- n limiting case.
642 However, we have provided evidence that these solutions continue to be more efficient relative
643 to DbS for small values of n , and more importantly, they remain nearly optimal even at relatively
644 low values of n , in the range of values that might be relevant to explain human experimental data
645 (*Appendix 7*).

646 Another key implication of our results is that we provide an alternative explanation to the usual
647 conception of noise as the main cause of behavioral performance degradation, where noise is usu-
648 ally artificially added to models of decision behavior to generate the desired variability (*Ratcliff*
649 *and Rouder, 1998; Wang, 2002*). On the contrary, our work makes it formally explicit why a system
650 that evolved to encode information based on binary codes must be necessarily noisy, also reveal-
651 ing how the system could take advantage of its unavoidable noisy properties (*Faisal et al., 2008*)
652 to optimize decision behavior (*Tsetsos et al., 2016*). Here it is important to highlight that this con-
653 clusion is drawn from a purely homogeneous neural circuit – in other words, a circuit in which all
654 neurons have the same properties (in our case, the same threshold activation thresholds). This
655 is not what is typically observed, as neural circuits are typically very heterogeneous. However, in
656 the neural circuit that we consider here, it could mean that the firing thresholds can vary across
657 neurons (*Orbán et al., 2016*), which could be used by the system to optimize the required variabil-
658 ity of binary neural codes. Interestingly, it has been shown in recent work that stochastic discrete
659 events also serve to optimize information transmission in neural population coding (*Ashida and*
660 *Kubo, 2010; Nikitin et al., 2009; Schmerl and McDonnell, 2013*). Crucially, in our work we provide
661 a direct link of the necessity of noise for systems that aim at optimizing decision behavior under
662 our encoding and limited-capacity assumptions, which can be seen as algorithmic specifications
663 of the more realistic population coding specifications mentioned above (*Nikitin et al., 2009*). We
664 acknowledge that based on the results of our work, we cannot confirm whether this is the case for
665 higher order neural circuits, however, we leave it as an interesting theoretical formulation, which
666 could be addressed in future work.

667 Interestingly, our results could provide an alternative explanation of the recent controversial
668 finding that dynamics of a large proportion of LIP neurons likely reflect binary (discrete) coding
669 states to guide decision behavior (*Latimer et al., 2015; Zoltowski et al., 2019*). Based on this poten-
670 tial link between our and their work, our theoretical framework generates testable predictions that
671 could be investigated in future neurophysiological work. For instance, noise distribution in neural
672 circuits should dynamically adapt according to the prior distribution of inputs and goals of the or-
673 ganism. Consequently, the rate of “step-like” coding in single neurons should also be dynamically
674 adjusted (perhaps optimally) to statistical regularities and behavioral goals.

675 Our results are closely related to Decision by Sampling (DbS), which is an influential account
676 of decision behavior derived from principles of retrieval and memory comparison by taking into
677 account the regularities of the environment, and also encodes information based on binary codes
678 (*Stewart et al., 2006*). We show that DbS represents a special case of our more general efficient
679 sampling framework, that uses a rule that is similar to (though not exactly like) the optimal en-
680 coding rule that assumes full (or costless) adaptation to the prior statistics of the environment.
681 In particular, we show that DbS might well be the most efficient sampling algorithm, given that
682 a reduction in the full representation of the prior distribution might not come at a great loss in
683 performance. Interestingly, our experimental results (discussed in more detail below) also provide
684 support for the hypothesis that numerosity perception is efficient in this particular way. Crucially,
685 DbS automatically adjusts the encoding in response to changes in the frequency distribution from
686 which exemplars are drawn in approximately the right way, while providing a simple answer to the
687 question of how such adaptation of the encoding rule to a changing frequency distribution occurs,
688 at a relatively low cost.

689 On a related line of work, *Bhui and Gershman (2018)* develop a similar, but different specifica-
690 tion of DbS, in which they also consider only a finite number of samples that can be drawn from
691 the prior distribution to generate a percept, and ask what kind of algorithm would be required to
692 improve coding efficiency. However, their implementation differs from ours in various important
693 ways (see *Appendix 8* for a detailed discussion). One of the main distinctions is that they consider
694 the case in which only a finite number of samples can be drawn from the prior and show that a vari-
695 ant of DbS with kernel-smoothing is superior to its standard version. However, a key difference to
696 our implementation is that they allow the kernel-smoothed quantity (computed by comparing the

697 input v with a sample \tilde{v} from the prior distribution) to vary continuously between 0 and 1, rather
698 than having to be either 0 or 1 as in our implementation (**Figure 1**). Thus, they show that coding
699 efficiency can be improved by allowing a more flexible implementation of the coding scheme for
700 the case when the agent is allowed to draw few samples from the prior distribution (**Appendix 8**).
701 On the other hand, we restrict our framework to a coding scheme that is only allowed to encode
702 information based on zeros or ones, where we show that coding efficiency can be improved rel-
703 ative to DbS only under a more complete knowledge of the prior distribution, where the optimal
704 solutions can be formally derived in the large- n limit. Nevertheless, we have shown that even un-
705 der the operation of few sampling units, the optimal rules will be still superior to the standard DbS
706 (if the agent has fully adapted to the statistics of the environment in a given context), even when a
707 few number of processing units are available to generate decision relevant percepts.

708 We tested these resource-limited coding frameworks in non-symbolic numerosity discrimina-
709 tion, a fundamental cognitive function for behavior in humans and other animals, which may have
710 emerged during evolution to support fitness maximization (**Nieder, 2020**). Here, we find that the
711 way in which the precision of numerosity discrimination varies with the size of the numbers being
712 compared is consistent with the hypothesis that the internal representations on the basis of which
713 comparisons are made are sample-based. In particular, we find that the encoding rule varies de-
714 pending on the frequency distribution of values encountered in a given environment, and that this
715 adaptation occurs fairly quickly once the frequency distribution changes.

716 This adaptive character of the encoding rule differs, for example, from the common hypothe-
717 sis of a logarithmic encoding rule (independent of context), which we show fits our data less well.
718 Nonetheless, we can reject the hypothesis of full optimality of the encoding rule for each distribu-
719 tion of values used in our experiments, even after subjects have had extensive experience with a
720 given distribution. Thus, a possible explanation of why DbS is the favored model in our numerosity
721 task is that accuracy and reward maximization requires optimal adaptation of the noise distribution
722 based on our imposed prior, requiring complex neuroplastic changes to be implemented, which
723 are in turn metabolically costly (**Buchanan et al., 2013**). Relying on samples from memory might
724 be less metabolically costly as these systems are plastic in short time scales, and therefore a rel-
725 atively simpler heuristic to implement allowing more efficient adaptation. Here it is important to
726 emphasize, as it has been discussed in the past (**Tajima et al., 2016; Polanía et al., 2015**), that for
727 decision-making systems beyond the perceptual domain, the identity of the samples is unclear. We
728 hypothesize, that information samples derive from the interaction of memory on current sensory
729 evidence depending on the retrieval of relevant samples to make predictions about the outcome
730 of each option for a given behavioral goal (therefore also depending on the encoding rule that
731 optimizes a given behavioral goal).

732 Interestingly, it was recently shown that in a reward learning task, a model that estimates val-
733 ues based on memory samples from recent past experiences can explain the data better than
734 canonical incremental learning models (**Bornstein et al., 2017**). Based on their and our findings,
735 we conclude that sampling from memory is an efficient mechanism for guiding choice behavior, as
736 it allows quick learning and generalization of environmental contexts based on recent experience
737 without significantly sacrificing behavioral performance. However, it should be noted that relying
738 on such mechanisms alone might be suboptimal from a performance- and goal-based point of
739 view, where neural calibration of optimal strategies may require extensive experience, possibly
740 via direct interactions between sensory, memory and reward systems (**Gluth et al., 2015; Saleem
741 et al., 2018**).

742 Taken together, our findings emphasize the need of studying optimal models, which serve as
743 anchors to understand the brain's computational goals without ignoring the fact that biological
744 systems are limited in their capacity to process information. We addressed this by proposing a
745 computational problem, elaborating an algorithmic solution, and proposing a minimalistic imple-
746 mentational architecture that solves the resource-constrained problem. This is essential, as it helps
747 to establish frameworks that allow comparing behavior not only across different tasks and goals,

748 but also across different levels of description, for instance, from single cell operation to observed
749 behavior (*Marr, 1982*). We argue that this approach is fundamental to provide benchmarks for
750 human performance that can lead to the discovery of alternative heuristics (*Qamar et al., 2013*;
751 *Gardner, 2019*) that could appear to be in principle suboptimal, but that might be in turn the opti-
752 mal strategy to implement if one considers cognitive limitations and costs of optimal adaptation.
753 We conclude that the understanding of brain function and behavior under a principled research
754 agenda, which takes into account decision mechanisms that are biologically feasible, will be essen-
755 tial to accelerate the elucidation of the mechanisms underlying human cognition.

756 **Methods and Materials**

757 **Participants**

758 The study tested young healthy volunteers with normal or corrected-to-normal vision (total n=20,
759 age 19-36 years, 9 females: n=7 in experiment 1, 2 females; n=6 new participants in experiment 2, 3
760 females; n=7 new participants in experiment 3, 4 females). Participants were randomly assigned to
761 each experiment and no participant was excluded from the analyses. Participants were instructed
762 about all aspects of the experiment and gave written informed consent. None of the participants
763 suffered from any neurological or psychological disorder or took medication that interfered with
764 participation in our study. Participants received monetary compensation for their participation
765 in the experiment partially related to behavioral performance (see below). The experiments con-
766 formed to the Declaration of Helsinki and the experimental protocol was approved by the Ethics
767 Committee of the Canton of Zurich (BASEC: 2018-00659).

768 **Experiment 1**

769 Participants (n=7) carried out a numerosity discrimination task for four consecutive days for ap-
770 proximately one hour per day. Each daily session consisted of a training run followed by 8 runs of
771 75 trials each. Thus, each participant completed ~2,400 trials across the four days of experiment.

772 After a fixation period (1-1.5s jittered), two clouds of dots (left and right) were presented on the
773 screen for 200ms. Participants were asked to indicate the side of the screen where they perceived
774 more dots. Their response was kept on the screen for 1 second followed by feedback consisting of
775 the symbolic number of dots in each cloud as well as the monetary gains and opportunity losses of
776 the trial depending on the experimental condition. In the *value* condition, participants were explic-
777 itly informed that each dot in a cloud of dots corresponded to 1 Swiss Franc (CHF). Participants were
778 informed that they would receive the amount in CHF corresponding to the total number of dots on
779 the chosen side. At the end of the experiment a random trial was selected and they received the
780 corresponding amount. In the *accuracy* condition, participants were explicitly informed that they
781 could receive a fixed reward (15 Swiss Francs (CHF)) for each correct trial. This fixed amount was
782 selected such that it approximately matched the expected reward received in the value condition
783 (as tested in pilot experiments). At the end of the experiment, a random trial was selected and
784 they would receive this fixed amount if they chose the cloud with more dots (i.e. the correct side).
785 Each condition lasted for two consecutive days with the starting condition randomized across par-
786 ticipants. Only after completing all four experiment days, participants were compensated for their
787 time with 20 CHF per hour, in addition to the money obtained based on their decisions on each
788 experimental day.

789 **Experiment 2**

790 Participants (n=6) carried out a numerosity discrimination task where each of four daily sessions
791 consisted of 16 runs of 40 trials each, thus each participant completed ~2,560 trials. A key differ-
792 ence with respect to Experiment 1 is that participants had to accumulate points based on their
793 decisions and had to reach a predetermined threshold on each run. The rules of point accumula-
794 tion depended on the experimental condition. In the perceptual condition, a fixed amount of points
795 was awarded if the participants chose the cloud with more dots. In this condition, participants were

796 instructed to accumulate a number of points and reach a threshold given a limited number of tri-
797 als. Based on the results obtained in Experiment 1, the threshold corresponded to 85% of correct
798 trials in a given run, however the participants were unaware of this. If the participants reached this
799 threshold, they were eligible for a fixed reward (20 CHF) as described in Experiment 1. In the value
800 condition, the number of points received was equal to the number of dots in the cloud, however,
801 contrary to experiment 1, points were only awarded if the participant chose the cloud with the
802 most dots. Participants had to reach a threshold that was matched in the expected collection of
803 points of the perceptual condition. As in Experiment 1, each condition lasted for two consecutive
804 days with the starting condition randomized across participants. Only after completing all the four
805 days of the experiment, participants were compensated for their time with 20 CHF per hour, in
806 addition to the money obtained based on their decisions on each experimental day.

807 Experiment 3

808 The design of Experiment 3 was similar to the value condition of Experiment 2 (n=7 participants)
809 and was carried out over three consecutive days. The key difference between Experiment 3 and
810 Experiments 1-2 was the shape of the prior distribution $f(v)$ that was used to draw the number of
811 dots for each cloud in each trial (see below).

812 Stimuli statistics and trial selection

813 For all experiments, we used the following parametric form of the prior distribution

$$f(v) = c(1 - v)^\alpha, \quad (10)$$

814 initially defined in the interval [0,1] for mathematical tractability in the analytical solution of the
815 encoding rules $\theta(v)$ (see below), with $\alpha > 0$ determining the shape of the distribution, and c is a
816 normalizing constant. For Experiments 1 and 2 the shape parameter was set to $\alpha = 2$, and for
817 Experiment 3 was set to $\alpha = 1$. i.i.d. samples drawn from this distribution were then multiplied by
818 50, added an offset of 5, and finally were rounded to the closest integer (i.e., the numerosity values
819 in our experiment ranged from $v_{\min} = 5$ to $v_{\max} = 55$). The pairs of dots on each trial were determined
820 by sampling from a uniform density window in the CDF space (Eq. 10 is its corresponding PDF).
821 The pairs of dots in each trial were selected with the conditions that, first, their distance in the CDF
822 space was less than a constant (0.25, 0.28 and 0.23 for Experiments 1, 2 and 3 respectively), and
823 second, the number of dots in both clouds was different. **Figure 3c** illustrates the probability that
824 a pair of choice alternatives was selected for a given trial in Experiments 1 and 2.

825 Power analyses and model recovery

826 Given that adaptation dynamics in sensory systems often require long-term experience with novel
827 prior distributions, we opted for maximizing the number of trials for a relatively small number of
828 participants per experiment, as it is commonly done for this type of psychophysical experiments
829 (*Brunton et al., 2013; Stocker and Simoncelli, 2006; Zylberberg et al., 2018*). Note that based on the
830 power analyses described below, we collected in total ~45,000 trials across the three Experiments,
831 which is above the average number of trials typically collected in human studies.

832 In order to maximize statistical power in the differentiation of the competing encoding rules,
833 we generated 10,000 sets of experimental trials for each encoding rule and selected the sets of
834 trials with the highest discrimination power (i.e. largest differences in Log-Likelihood) between the
835 encoding models. In these power analyses, we also investigated what was the minimum number
836 of trials that would allow accurate generative model selection at the individual level. We found
837 that ~1,000 trials per participant in each experimental condition would be sufficient to predict ac-
838 curately ($P > 0.95$) the true generative model. Based on these analyses, we decided to collect at least
839 1,200 trials per participant and condition (perceptual and value) in each of the three experiments.
840 Model recovery analyses presented in **Figure 3d** illustrate the result of our power analyses (see
841 also **Figure 3–Figure Supplement 1**).

842 Apparatus

843 Eyetracking (EyeLink 1000 Plus) was used to check the participants fixation during stimulus presen-
 844 tation. When participants blinked or move their gaze (more than 2° of visual angle) away from the
 845 fixation cross during the stimulus presentation the trial was canceled (only 212 out of 45,600 trials
 846 were canceled, i.e., < 0.5% of the trials). Participants were informed when a trial was canceled and
 847 were encouraged not to do so as they would not receive any reward for this trial. A chinrest was
 848 used to keep the distance between the participants and the screen constant (55cm). The task was
 849 run using Psychtoolbox Version 3.0.14 on Matlab 2018a. The diameter of the dots varied between
 850 0.42° and 1.45° of visual angle. The center of each cloud was positioned 12.6° of visual angle hor-
 851 izontally from the fixation cross and had a maximum diameter of 19.6° of visual angle. Following
 852 previous numerosity experiments (*Berg et al., 2017; Izard and Dehaene, 2008*), either the average
 853 dot size or the total area covered by the dots was maintained constant in both clouds for each
 854 trial. The color of each dot (white or black) was randomly selected for each dot. Stimuli set were
 855 different for each participant but identical between the two conditions.

856 Encoding rules and model fits

857 The parametrization of the prior $f(v)$ (Eq. 10) allows tractable analytical solutions of the encoding
 858 rules $\theta_A(v)$, $\theta_R(v)$ and $\theta_D(v)$, that correspond to Accuracy maximization, Reward maximization, and
 859 DbS, respectively:

$$\theta_A(v) = \sin \left[\frac{\pi}{2} (1 - (1 - v)^{\alpha+1}) \right]^2 \quad (11)$$

$$\theta_R(v) = \sin \left[\frac{\pi}{2} (1 + (v - 1)((1 - v)^\alpha)^{2/3}) \right]^2 \quad (12)$$

$$\theta_D(v) = 1 - (1 - v)^{\alpha+1} \quad (13)$$

862 Graphical representation of the respective encoding rules is shown in **Figure 3e** for Experiments
 863 1 and 2. Given an encoding rule $\theta(v)$, we now define the decision rule. The goal of the decision
 864 maker in our task is always to decide which of two input values v_1 and v_2 is larger. Therefore, the
 865 agent choses v_1 if and only if the internal readings $k_1 > k_2$. Following the definitions of expected
 866 value and variance of binomial variables, and approximating for large n (see **Appendix 2**), the prob-
 867 ability of choosing v_1 is given by

$$P_{\text{choose } v_1} \approx \Phi \left(\frac{\theta_1 - \theta_2}{\sqrt{\frac{\theta_1(1-\theta_1) + \theta_2(1-\theta_2)}{n}}} \right) \quad (14)$$

868 where $\Phi()$ is the standard CDF, and θ_1 and θ_2 are the encoding rules for the input values v_1 and v_2 ,
 869 respectively. Thus, the choice structure is the same for all models, only differing in their encoding
 870 rule. The three models generate different qualitative performance predictions for a given number
 871 of samples n (**Figure 3f**).

872 Crucially, this probability decision rule (Eq. 14) can be parsimoniously extended to include po-
 873 tential side biases independent of the encoding process as follows

$$P_{\text{choose } v_1} \approx \Phi \left(\frac{\theta_1 - \theta_2}{\sqrt{\frac{\theta_1(1-\theta_1) + \theta_2(1-\theta_2)}{n}}} + \beta_0 \right) \quad (15)$$

874 where β_0 is the bias term. This is the base model used in our work. We were also interested in
 875 studying whether choice history effects (*Abrahamyan et al., 2016; Talluri et al., 2018*) may have in-
 876 fluence in our task, thus possibly affecting the conclusions that can be drawn from the base model.

877 Therefore, we extended this model to incorporate the effect of decision learning and choices from
878 the previous trial

$$P_{\text{choose } v_1} \approx \Phi \left(\frac{\theta_1 - \theta_2}{\sqrt{\frac{\theta_1(1-\theta_1) + \theta_2(1-\theta_2)}{n}}} + \beta_0 + \beta^L a_{t-1} r_{t-1} + \beta^{\text{Ch}} a_{t-1} \right), \quad (16)$$

879 where a_{t-1} is the choice made on the previous trial (+1 for left choice and -1 for right choice) and
880 r_{t-1} is the “outcome learning” on the previous trial (+1 for correct choice and -1 for incorrect choice).
881 β^L and β^{Ch} capture the effect of decision learning and choice in the previous trial, respectively.

882 Given that the choice structure is the same for all three sampling models considered here, we
883 can naturally address the question of what decision rule the participants favor via a latent-mixture
884 model. We implemented this model based on a hierarchical Bayesian modelling (HBM) approach.
885 The base-rate probabilities for the three different encoding rules at the population level are repre-
886 sented by the vector $\boldsymbol{\pi}$, so that π_m is the probability of selecting encoding rule model m . We initialize
887 the model with an uninformative prior given by

$$\boldsymbol{\pi} \sim \text{Dirichlet}(1_{m=1}, 1_{m=2}, 1_{m=3}).$$

888 This base-rate is updated based on the empirical data, where we allow each participant s to draw
889 from each model categorically based on the updated base-rate

$$m_s \sim \text{Categorical}(\boldsymbol{\pi}),$$

890 where the encoding rule θ for model m is given by

$$\theta_{m,s} = \begin{cases} \theta_A, & m = 1 \\ \theta_R, & m = 2 \\ \theta_D, & m = 3 \end{cases}$$

891 The selected rule was then fed into equations 15 or 16 to determine the probability of selecting a
892 cloud of dots. The number of samples n was also estimated within the same HBM with population
893 mean μ and standard deviation σ initialized based on uninformative priors with plausible ranges

$$\begin{aligned} \mu_n &\sim \text{Uniform}(1, 1000) \\ \sigma_n &\sim \text{Uniform}(0.01, 1000) \end{aligned}$$

allowing each participant s to draw from this population prior assuming that n is normally dis-
tributed at the population level

$$n_s \sim \text{Normal}(\mu_n, \sigma_n)$$

894 Similarly, the latent variables β in equations 15 and 16 were estimated by setting popu-
895 lation mean μ_β and standard deviation σ_β initialized based on uninformative priors

$$\begin{aligned} \mu_\beta &\sim \text{Uniform}(-10, 10) \\ \sigma_\beta &\sim \text{Uniform}(0.01, 100) \end{aligned}$$

allowing each participant s to draw from this population prior assuming that β is normally dis-
tributed at the population level

$$\beta_s \sim \text{Normal}(\mu_\beta, \sigma_\beta)$$

896 In all the results reported in **Figure 3** and **Figure 4**, the value of the shape parameter of the
897 prior was set to its true value $\alpha = 2$. The estimation of α in **Figure 5a** was investigated with a similar
898 hierarchical approach, allowing each participant to sample from the normal population distribution
899 with uninformative priors over the population mean and standard deviation

$$\begin{aligned} \mu_\alpha &\sim \text{Uniform}(0.01, 20) \\ \sigma_\alpha &\sim \text{Uniform}(0.0001, 100) \end{aligned}$$

900 The choice rule of the standard logarithmic model of numerosity discrimination is given by

$$P_{\text{choose } v_1} = \Phi\left(\frac{\log(v_1) - \log(v_2)}{\sigma\sqrt{2}}\right), \quad (17)$$

where σ is the internal noise in the logarithmic space. This model was extended to incorporate bias and choice history effects in the same way as implemented in the sampling models. Here we emphasize that all sampling and log models have the same degrees of freedom, where performance is mainly determined by n in the sampling models and Weber's fraction σ in the log model, and biases are determined by parameters β . For all above-mentioned models, the trial-by-trial likelihood of the observed choice (i.e. the data) given probability of a decision was based on a Bernoulli process

$$y_{t,s} \sim \text{Bernoulli}(P_{\text{choose } v_1})$$

901 where $y_{t,s} \in \{0, 1\}$ is the decision of each participant s in each trial t . In order to allow for prior
902 adaptation, the model fits presented in **Figure 3** and **Figure 4** were fit starting after a fourth of the
903 daily trials (corresponding to 150 trials for experiment 1 and 160 trials for experiment 2) to allow
904 for prior adaptation and fixing the shape parameter to its true generative value $\alpha = 2$.

905 The dynamics of adaptation (**Figure 5**) were studied by allowing the shape parameter α to evolve
906 through trial experience using all trials collected on each experiment day. This was studied using
907 the following function

$$\alpha_t = \delta + \eta e^{-t/\tau}, \quad (18)$$

908 where δ represents a possible target adaptation value of α , t is the trial number, and η , τ determine
909 the shape of the adaptation. Therefore, the encoding rule of the DbS model also changed trial-to-
910 trial

$$\theta_D^t(v) = 1 - (1 - v)^{\alpha_t + 1}. \quad (19)$$

911 Adaptation was tested based on the hypothesis that participants initially use a logarithmic dis-
912 crimination rule (Eq. 17) (this strategy also allowed improving identification of the adaptation dy-
913 namics). Therefore, Eq. 18 was parametrized such that the initial value of the shape parameter
914 ($\alpha_{t=0}$) guaranteed that discriminability between the DbS and the logarithmic rule was as close as
915 possible. This was achieved by finding the value of α in the DbS encoding rule (θ_D) that minimizes
916 the following expression

$$\sum_{t=1}^T \left[\left(\frac{\theta_D(v_{1,t}) - \theta_D(v_{2,t})}{\sqrt{\theta_D(v_{1,t})(1 - \theta_D(v_{1,t})) + \theta_D(v_{2,t})(1 - \theta_D(v_{2,t}))}} \right) - (\log(v_{1,t}) - \log(v_{2,t})) \right]^2, \quad (20)$$

917 where $v_{1,t}$ and $v_{2,t}$ are the numerosity inputs for each trial t . This expression was minimized based
918 on all trials generated in Experiments 1-3 (note that minimizing this expression does not require
919 knowledge of the sensitivity levels σ and n for the log and DbS models, respectively). We found that
920 the shape parameter value that minimizes Eq. 20 is $\alpha = 2.58$. Based on our prior $f(v)$ parametriza-
921 tion (Eq. 10), this suggests that the initial prior is more skewed than the priors used in Experiments
922 1-3 (**Figure 5b**). This is an expected result given that log-normal priors – typically assumed in nu-
923 merosity tasks – are also highly skewed. We fitted the δ parameter independently for Experiments
924 1-2 and Experiments 3 but kept the τ parameter shared across all experiments. If adaptation is tak-
925 ing place, we hypothesized that the asymptotic value δ of the shape parameter α should be larger
926 for Experiments 1-2 compared to Experiment 3.

927 Posterior inference of the parameters in all the hierarchical models described above was per-
928 formed via the Gibbs sampler using the Markov Chain Monte Carlo (MCMC) technique implemented
929 in JAGS. For each model, a total of 50,000 samples were drawn from an initial burn-in step and sub-
930 sequently a total of new 50,000 samples were drawn for each of three chains (samples for each
931 chain were generated based on a different random number generator engine, and each with a dif-
932 ferent seed). We applied a thinning of 50 to this final sample, thus resulting in a final set of 1,000

933 samples for each chain (for a total of 3,000 pooling all 3 chains). We conducted Gelman–Rubin tests
934 for each parameter to confirm convergence of the chains. All latent variables in our Bayesian mod-
935 els had $\hat{R} < 1.05$, which suggests that all three chains converged to a target posterior distribution.
936 We checked via visual inspection that the posterior population level distributions of the final MCMC
937 chains converged to our assumed parametrizations. When evaluating different models, we are in-
938 terested in the model's predictive accuracy for unobserved data, thus it is important to choose a
939 metric for model comparison that considers this predictive aspect. Therefore, in order to perform
940 model comparison, we used a method for approximating leave-one-out cross-validation (LOO) that
941 uses samples from the full posterior (*Vehtari et al., 2016*). These analyses were repeated using an
942 alternative Bayesian metric: the WAIC (*Vehtari et al., 2016*).

943 **Data availability**

944 Data and code that support the findings of this study will be made available via an open repository.

945 **Acknowledgments**

946 This work was supported by an ERC starting grant (ENTRAINER) to R.P and by a grant of the U.S. Na-
947 tional Science Foundation to M.W. This project has received funding from the European Research
948 Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant
949 agreement No. 758604).

950 **Competing Interests Statement**

951 The authors declare no competing financial interests.

952 **References**

- 953 **Abrahamyan A**, Silva LL, Dakin SC, Carandini M, Gardner JL. Adaptable history biases in human perceptual
954 decisions. *Proc Natl Acad Sci.* 2016; 113:3548– 3557.
- 955 **Ashida G**, Kubo M. Suprathreshold stochastic resonance induced by ion channel fluctuation. *Phys D Nonlinear*
956 *Phenom.* 2010; 239:327–334.
- 957 **Attneave F.** Some informational aspects of visual perception. *Psychol Rev.* 1954; 61:183.
- 958 **Berg R**, Lindskog M, Poom L, Winman A. Recent is more: A negative time-order effect in nonsymbolic numerical
959 judgment. *J Exp Psychol Hum Percept Perform.* 2017; 43:1084–1097.
- 960 **Bhui R**, Gershman SJ. Decision by sampling implements efficient coding of psychoeconomic functions. *Psychol*
961 *Rev.* 2018; 125:985–1001.
- 962 **Bornstein AM**, Khaw MW, Shohamy D, Daw ND. Reminders of past choices bias decisions for reward in humans.
963 *Nat Commun.* 2017; 8:15958.
- 964 **Brunton BW**, Botvinick MM, Brody CD. Rats and Humans Can Optimally Accumulate Evidence for Decision-
965 Making. *Science.* 2013; 80-.). 340:95–98.
- 966 **Brus J**, Heng JA, Polanía R. Weber's Law: A Mechanistic Foundation after Two Centuries. *Trends in Cognitive*
967 *Sciences.* 2019; 23(11):906–908.
- 968 **Buchanan KL**, Grindstaff JL, Pravosudov VV. Condition dependence, developmental plasticity, and cognition:
969 implications for ecology and evolution. *Trends Ecol Evol.* 2013; 28:290–296.
- 970 **Butterworth B**, Gallistel CR, Vallortigara G. Introduction: The origins of numerical abilities. *Philos Trans R Soc*
971 *B Biol Sci.* 2018; 373:20160507.
- 972 **Clarke BS**, Barron AR. Jeffreys' prior is asymptotically least favorable under entropy risk. *J Stat Plan Inference.*
973 1994; 41:37–60.
- 974 **Dehaene S.** The neural basis of the Weber-Fechner law: a logarithmic mental number line. *Trends Cogn Sci.*
975 2003; 7:145–147.

- 976 **Dutilh G**, Rieskamp J. Comparing perceptual and preferential decision making. *Psychon Bull Rev.* 2016;
977 23:723–737.
- 978 **Faisal AA**, Selen LPJ, Wolpert DM. Noise in the nervous system. *Nat Rev Neurosci.* 2008; 9:292–303.
- 979 **Ganguli D**, Simoncelli EP. Efficient sensory encoding and Bayesian inference with heterogeneous neural pop-
980 ulations. *Neural Comput.* 2014; 26:2103–2134.
- 981 **Ganguli D**, Simoncelli EP, Neural and perceptual signatures of efficient sensory coding; 2016.
- 982 **Gardner JL**. Optimality and heuristics in perceptual neuroscience. *Nat Neurosci.* 2019; 1.
- 983 **Gluth S**, Sommer T, Rieskamp J, Büchel C. Effective Connectivity between Hippocampus and Ventromedial
984 Prefrontal Cortex Controls Preferential Choices from Memory. *Neuron.* 2015; 86:1078–1090.
- 985 **Grueschow M**, Polania R, Hare TA, Ruff CC. Automatic versus Choice-Dependent Value Representations in the
986 Human Brain. *Neuron.* 2015; 85:874–885.
- 987 **Hawkes K**, O’Connell JF, Jones NG, Alvarez H, Charnov EL, Hof PR, Wildman DE, Sherwood CC, Leonard WR,
988 Lange N. Grandmothering, menopause, and the evolution of human life histories. *Proc Natl Acad Sci U S A.*
989 1998; 95:1336–1339.
- 990 **Izard V**, Dehaene S. Calibrating the mental number line. *Cognition.* 2008; 106:1221–1247.
- 991 **Keung W**, Hagen TA, Wilson RC. Regulation of evidence accumulation by pupil-linked arousal processes. *Nat*
992 *Hum Behav.* 2019; 1.
- 993 **Khaw MW**, Li Z, Woodford M. Review of Economic Studies. In Press; .
- 994 **Kolling N**, Wittmann M, Rushworth MFSFS. Multiple Neural Mechanisms of Decision Making and Their Compe-
995 tition under Changing Risk Pressure. *Neuron.* 2014; 81:1190–1202.
- 996 **Latimer KW**, Yates JL, Meister MLR, Huk AC, Pillow JW. Single-trial spike trains in parietal cortex reveal discrete
997 steps during decision-making. *Science.* 2015; 349:184–187.
- 998 **Laughlin S**. A Simple Coding Procedure Enhances a Neuron’s Information Capacity. *Zeitschrift Für Naturforsch*
999 *C.* 1981; 36:910–912.
- 1000 **Louie K**, Glimcher PW. Efficient coding and the neural representation of value. *Ann N Y Acad Sci.* 2012;
1001 1251:13–32.
- 1002 **Marr D**, *Vision: A computational investigation into the human representation and processing of visual infor-*
1003 *mation*; 1982.
- 1004 **McDonnell MD**, Stocks NG, Abbott D. Optimal stimulus and noise distributions for information transmission
1005 via suprathreshold stochastic resonance. *Phys Rev E.* 2007; 75:061105.
- 1006 **Młynarski W**, Hermundstad AM. Adaptability and efficiency in neural coding. In: *BioRxiv 669200*; 2019.
- 1007 **Navarrete A**, Schaik CP, Isler K. Energetics and the evolution of human brain size. *Nature.* 2011; 480:91–93.
- 1008 **Nieder A**. The Adaptive Value of Numerical Competence. *Trends Ecol Evol.* 2020; 0.
- 1009 **Nieder A**, Dehaene S. Representation of number in the brain. *Annu Rev Neurosci.* 2009; 32:185–208.
- 1010 **Nieder A**, Miller EK. Coding of cognitive magnitude: compressed scaling of numerical information in the pri-
1011 mate prefrontal cortex. *Neuron.* 2003; 37:149–157.
- 1012 **Nikitin AP**, Stocks NG, Morse RP, McDonnell MD. Neural Population Coding Is Optimized by Discrete Tuning
1013 Curves. *Phys Rev Lett.* 2009; 103:138101.
- 1014 **Niven JE**, Laughlin SB. Energy limitation as a selective pressure on the evolution of sensory systems. *J Exp Biol.*
1015 2008; 211:1792–1804.
- 1016 **Norman DA**. Toward a theory of memory and attention. *Psychol Rev.* 1968; 75:522–536.
- 1017 **Orbán G**, Berkes P, Fiser J, Lengyel M. Neural Variability and Sampling-Based Probabilistic Representations in
1018 the Visual Cortex. *Neuron.* 2016; 92:530–543.

- 1019 **Pardo-Vazquez JL**, Castiñeiras-de Saa JR, Valente M, Damião I, Costa T, Vicente MI, Mendonça AG, Mainen ZF,
1020 Renart A. The mechanistic foundation of Weber's law. *Nature neuroscience*. 2019; p. 1–10.
- 1021 **Park IM**, Pillow JW, Bayesian Efficient Coding; 2017.
- 1022 **Piazza M**, Pinel P, Le Bihan D, Dehaene S. A magnitude code common to numerosities and number symbols
1023 in human intraparietal cortex. *Neuron*. 2007; 53:293–305.
- 1024 **Pirrone A**, Stafford T, Marshall JAR. When natural selection should optimize speed-accuracy trade-offs. *Front*
1025 *Neurosci*. 2014; 8:73.
- 1026 **Polanía R**, Krajbich I, Grueschow M, Ruff CC. Neural oscillations and synchronization differentially support
1027 evidence accumulation in perceptual and value-based decision-making. *Neuron*. 2014; 82:709–720.
- 1028 **Polanía R**, Moisa M, Opitz A, Grueschow M, Ruff CC. The precision of value-based choices depends causally on
1029 fronto-parietal phase coupling. *Nat Commun*. 2015; 6:8090.
- 1030 **Polanía R**, Woodford M, Ruff CC. Efficient coding of subjective value. *Nat Neurosci*. 2019; 22:134–142.
- 1031 **Qamar AT**, Cotton RJ, George RG, Beck JM, Prezhdo E, Laudano A, Tolia AS, Ma WJ. Trial-to-trial,
1032 uncertainty-based adjustment of decision boundaries in visual categorization. *Proc Natl Acad Sci*. 2013;
1033 110:20332–20337.
- 1034 **Ratcliff R**, Rouder JN. Modeling Response Times for Two-Choice Decisions. *Psychol Sci*. 1998; 9:347–356.
- 1035 **Rustichini A**, Conen KE, Cai X, Padoa-Schioppa C. Optimal coding and neuronal adaptation in economic deci-
1036 sions. *Nat Commun*. 2017; 8:1208.
- 1037 **Saleem AB**, Diamanti EM, Fournier J, Harris KD, Carandini M. Coherent encoding of subjective spatial position
1038 in visual cortex and hippocampus. *Nature*. 2018; 562:124–127.
- 1039 **Schmerl BA**, McDonnell MD. Channel-noise-induced stochastic facilitation in an auditory brainstem neuron
1040 model. *Phys Rev E*. 2013; 88:052722.
- 1041 **Schreiber S**, Machens CK, Herz AVM, Laughlin SB. Energy-Efficient Coding with Discrete Stochastic Events.
1042 *Neural Comput*. 2002; 14:1323–1346.
- 1043 **Shadlen MNN**, Shohamy D. Decision Making and Sequential Sampling from Memory. *Neuron*. 2016;
1044 90:927–939.
- 1045 **Sharpee TO**. Optimizing Neural Information Capacity through Discretization. *Neuron*. 2017; 94:954–960.
- 1046 **Sharpee TO**, Calhoun AJ, Chalasani SH. Information theory of adaptation in neurons, behavior, and mood. *Curr*
1047 *Opin Neurobiol*. 2014; 25:47–53.
- 1048 **Stephan KE**, Penny WD, Daunizeau J, Moran RJ, Friston KJ. Bayesian model selection for group studies. *Neu-*
1049 *roimage*. 2009; 46:1004–1017.
- 1050 **Stewart N**, Chater N, Brown GDA. Decision by sampling. *Cogn Psychol*. 2006; 53:1–26.
- 1051 **Stewart N**, Reimers S, Harris AJL. On the Origin of Utility, Weighting, and Discounting Functions: How They Get
1052 Their Shapes and How to Change Their Shapes. *Manage Sci*. 2015; 61:687–705.
- 1053 **Stockler AA**, Simoncelli EP. Noise characteristics and prior expectations in human visual speed perception. *Nat*
1054 *Neurosci*. 2006; 9:578–585.
- 1055 **Stocks NG**, Allingham D, Morse RP. The application of suprathreshold stochastic resonance to cochlear implant
1056 coding. . 2002; .
- 1057 **Tajima S**, Drugowitsch J, Pouget A. Optimal policy for value-based decision-making. *Nat Commun*. 2016;
1058 7:12400.
- 1059 **Talluri BC**, Urai AE, Tsetsos K, Usher M, Donner TH. Confirmation Bias through Selective Overweighting of
1060 Choice-Consistent Evidence. *Curr Biol*. 2018; 28:3128–3135 8.
- 1061 **Tsetsos K**, Moran R, Moreland J, Chater N, Usher M, Summerfield C. Economic irrationality is optimal during
1062 noisy decision making. *Proc Natl Acad Sci*. 2016; 113:3102–3107.

- 1063 **Vehtari A**, Gelman A, Gabry J. Practical Bayesian model evaluation using leave-one-out cross-validation and
1064 WAIC. *Stat Comput.* 2016; p. 1–20.
- 1065 **Wang XJ**. Probabilistic decision making by slow reverberation in cortical circuits. *Neuron.* 2002; 36:955–968.
- 1066 **Weber EU**, Johnson EJ. Mindful judgment and decision making. *Annu Rev Psychol.* 2009; 60:53–85.
- 1067 **Wei XX**, Stocker AA. A Bayesian observer model constrained by efficient coding can explain ‘anti-
1068 Bayesian’ percepts. *Nat Neurosci.* 2015; 18:1509–1517.
- 1069 **Wei XX**, Stocker AA. Lawful relation between perceptual bias and discriminability. *Proc Natl Acad Sci U S A.*
1070 2017; 114:10244–10249.
- 1071 **Woodford M**. Modeling imprecision in perception, valuation and choice. *Annual Review of Economics.* In Press;
1072 .
- 1073 **Zoltowski DM**, Latimer KW, Yates JL, Huk AC, Pillow JW. Discrete Stepping and Nonlinear Ramping Dynamics
1074 Underlie Spiking Responses of LIP Neurons during Decision-Making. *Neuron.* 2019; 102:1249–1258 10.
- 1075 **Zylberberg A**, Wolpert DM, Shadlen MN. Counterfactual Reasoning Underlies the Learning of Priors in Decision
1076 Making. *Neuron.* 2018; 0.

Table 1. Resource parameter n fits.

Experiment	Condition	History effects	Model		
			$n_{Accuracy}$	n_{Reward}	n_{DbS}
1	V	not included	15.24 ± 3.09	17.54 ± 3.98	24.40 ± 5.16
2	V	not included	22.48 ± 2.43	27.58 ± 3.81	35.40 ± 3.44
1	P	not included	15.19 ± 3.99	17.84 ± 4.85	24.64 ± 6.59
2	P	not included	20.99 ± 1.59	24.22 ± 1.93	33.54 ± 2.45
1	P/V	not included	15.33 ± 3.41	17.25 ± 4.45	24.15 ± 5.75
2	P/V	not included	21.30 ± 0.96	25.27 ± 1.99	33.90 ± 1.51
1/2	V	not included	18.56 ± 2.04	22.05 ± 2.73	29.52 ± 3.25
1/2	P	not included	17.91 ± 2.09	20.66 ± 2.59	28.62 ± 3.51
1/2	P/V	not included	17.93 ± 1.87	21.03 ± 2.46	28.58 ± 3.04
1	V	included	15.50 ± 3.13	17.50 ± 3.91	24.68 ± 5.08
2	V	included	22.92 ± 2.37	28.07 ± 3.73	36.18 ± 2.91
1	P	included	15.41 ± 3.81	17.96 ± 4.88	24.70 ± 6.62
2	P	included	21.57 ± 1.71	24.88 ± 2.17	34.37 ± 2.93
1	P/V	included	15.16 ± 3.55	17.43 ± 4.39	24.30 ± 5.94
2	P/V	included	21.80 ± 0.92	25.81 ± 1.86	34.60 ± 1.40
1/2	V	included	18.86 ± 2.07	22.48 ± 2.75	29.85 ± 3.17
1/2	P	included	18.15 ± 2.17	21.11 ± 2.72	29.01 ± 3.47
1/2	P/V	included	18.22 ± 1.93	21.34 ± 2.50	29.12 ± 3.12

Fits of the resource parameter for the Accuracy, Reward and Decision by Sampling (DbS) models including data across experiments and conditions (Perceptual (P) or Value (V)) either including or ignoring choice history effects. The values represent the mean ± SD of the posterior distributions at the population level for parameter n . Note that Reward and in particular the DbS encoding models require a higher number of resources than the Accuracy model, which is coherent with the fact that the Accuracy model allocates its resources to maximize efficiency, therefore reducing the number of resources needed to reach a given accuracy. DbS has the highest values of n because it is the most inefficient model.

Table 1-source data 1.

1077 **Appendix 1**

1078

Infomax coding rule

We assume that the subjective perception of an environmental variable with value v is determined by n independent *samples* of a binary random variable, i.e. outcomes are either "high" (ones) or "low" (zeros) readings. Here, the probability θ of a "high" reading is the same on each draw, but can depend on the input stimulus value, via the function $\theta(v)$. Additionally, we assume that the input value v on a given trial is an independent draw from some prior distribution $f(v)$ in a given environment or context (with $F(v)$ being the corresponding cumulative distribution function). As we mentioned before, the choice of θ (i.e. encoding of the input value) depends on v . Now suppose that the mapping $\theta(v)$ (the encoding rule) is chosen so as to maximize the mutual information between the random variable v and the subjective value representation k . The mutual information is computed under the assumption that v is drawn from a particular prior distribution $f(v)$, and $\theta(v)$ is assumed to be optimized for this prior. The mutual information between v and k is defined as

$$I(v, k) = H(k) - H(k|v), \quad (21)$$

where the marginal entropy $H(k)$ quantifies the uncertainty of the marginal response distribution $P(k)$, and $H(k|v)$ is the average conditional entropy of k given v . The output distribution is given by

$$P(k) = \int_{v \in V} P(k|v) f(v) dv, \quad (22)$$

where $f(v)$ is defined as the input density function. For the encoding framework that we consider here which is given by the binomial channel, the conditional probability mass function of the output given the input is

$$P(k|v) = \binom{n}{k} \theta(v)^k (1 - \theta(v))^{n-k}, \quad k \in [0, 1, \dots, n]. \quad (23)$$

Thus, we have all the ingredients to write the expression of the mutual information

$$\begin{aligned} I(v, k) &= H(k) - H(k|v) \\ &= - \sum_{k=0}^n P(k) \log P(k) \\ &\quad - \left(- \int_{v \in V} f(v) \sum_{k=0}^n P(k|v) \log P(k|v) dv \right) \end{aligned} \quad (24)$$

We then seek to determine the encoding rule $\theta(v)$ that solves the optimization problem

$$\text{find } C = \max_{\{\theta(v)\}} I(v, k). \quad (25)$$

It can be shown that for large n , the mutual information between θ and k (hence the mutual information between v and k) is maximized if the prior distribution over θ is the Jeffreys prior (**Clarke and Barron, 1994**)

$$\text{Beta}(\theta; 0.5, 0.5) = \frac{1}{\pi \sqrt{\theta(1 - \theta)}}, \quad (26)$$

also known as the arcsine distribution. Hence, the mapping $\theta(v)$ induces a prior distribution over θ given by the arcsine distribution. This means that for each v , the encoding function

1118

1119

1120

1121

$\theta(v)$ must be such that

1122

1123

1124

1125

1126

1127

1128

1129

$$\begin{aligned} F(v) &= \int_0^{\theta(v)} \frac{1}{\pi \sqrt{\tilde{\theta}(1-\tilde{\theta})}} d\tilde{\theta} \\ &= \frac{2}{\pi} \arcsin(\sqrt{\theta(v)}). \end{aligned} \quad (27)$$

Solving for θ we finally obtain the optimal encoding rule

$$\theta(v) = \left[\sin\left(\frac{\pi}{2} F(v)\right) \right]^2. \quad (28)$$

1130 Appendix 2

1131 Accuracy maximization for a known prior distribution

1132 Here we derive the optimal encoding rule when the criterion to be maximized is the prob-
 1133 ability of a correct response in a binary comparison task, rather than mutual information
 1134 as in **Appendix 1**. As in **Appendix 1**, we assume that the prior distribution $f(x)$ from which
 1135 stimuli are drawn is known, and that the encoding rule is optimized for this particular dis-
 1136 tribution. (The case in which we wish the encoding rule to be robust to variations in the
 1137 distribution from which stimuli are drawn is instead considered in **Appendix 6**.) Note that
 1138 the objective assumed here corresponds to maximization of expected reward in the case
 1139 of a perceptual experiment in which a subject must indicate which of two presented magni-
 1140 tudes is greater, and is rewarded for the number of correct responses. (In **Appendix 5**, we
 1141 instead consider the encoding rule that would maximize expected reward if the subject's
 1142 reward is proportional to the magnitude selected by their response.)

1143 As above, we assume encoding by a binomial channel. The encoded value (number of
 1144 "high" readings) is given by k , which is consequently an integer between 0 and n . This is a
 1145 random variable with a binomial distribution with expected value and variance given by

$$1146 \quad \mathbb{E} \left[\frac{k}{n} | \theta \right] = \theta \quad \text{Var} \left[\frac{k}{n} | \theta \right] = \frac{\theta(1-\theta)}{n} \quad (29)$$

1147 Suppose that the task of the decision maker is to decide which of two input values v_1 and v_2
 1148 is larger. Assuming that v_1 and v_2 are encoded independently, then the decision maker chooses
 1149 v_1 if and only if the internal readings $k_1 > k_2$ (here we may suppose that the probability
 1150 of choosing stimulus 1 is 0.5 in the event that $k_1 = k_2$). Thus, the probability of choosing
 1151 stimulus 1 is:

$$1152 \quad \mathbb{P} \left(\frac{k_1}{n} > \frac{k_2}{n} | v_1, v_2 \right) + \frac{1}{2} \mathbb{P} \left(\frac{k_1}{n} = \frac{k_2}{n} | v_1, v_2 \right). \quad (30)$$

1153 In the case of large n , we can use a normal approximation to the binomial distribution to
 1154 obtain

$$1155 \quad \left(\frac{k_1}{n} - \frac{k_2}{n} \right) \sim \mathcal{N} \left(\theta_1 - \theta_2, \frac{\theta_1(1-\theta_1) + \theta_2(1-\theta_2)}{n} \right) \quad (31)$$

1156 and hence the probability of choosing v_1 is given by

$$1157 \quad \mathbb{P}_{\text{choose } v_1} \approx \Phi \left(\frac{\theta_1 - \theta_2}{\sqrt{\frac{\theta_1(1-\theta_1) + \theta_2(1-\theta_2)}{n}}} \right), \quad (32)$$

1158 where $\Phi(\cdot)$ is the standard CDF. Thus the probability of an incorrect choice (i.e. choosing the
 1159 item with the lower value) is approximately

$$1160 \quad \mathbb{P}_{\text{error}} \approx \Phi \left(-\frac{|\theta_1 - \theta_2|}{\sqrt{\frac{\theta_1(1-\theta_1) + \theta_2(1-\theta_2)}{n}}} \right) \quad (33)$$

1161 Now, suppose that the encoding rule, together with the prior distribution for v (the same for
 1162 both inputs, that are independent draws from the prior distribution) results in an ex-ante

distribution for θ (same for both goods) with density function $\hat{f}(\theta)$. Then the probability of error is given by

$$P_{\text{error}} \approx \int \int \Phi \left(-\frac{|\theta_1 - \theta_2|}{\sqrt{\frac{\theta_1(1-\theta_1) + \theta_2(1-\theta_2)}{n}}} \right) \hat{f}(\theta_1) \hat{f}(\theta_2) d\theta_1 d\theta_2 \quad (34)$$

Our goal is to evaluate Eq. 34 for any choice of the density $\hat{f}(\theta)$. First, we fix the value of θ_1 and integrate over θ_2 :

$$\begin{aligned} & \int_0^1 \Phi \left(-\frac{|\theta_1 - \theta_2|}{\sqrt{\theta_1(1-\theta_1) + \theta_2(1-\theta_2)}} \sqrt{n} \right) \hat{f}(\theta_2) d\theta_2 \\ &= \int_0^{\theta_1} \Phi \left(-\frac{\theta_2 - \theta_1}{\sqrt{\theta_1(1-\theta_1) + \theta_2(1-\theta_2)}} \sqrt{n} \right) \hat{f}(\theta_2) d\theta_2 \\ & \quad + \int_{\theta_1}^1 \Phi \left(-\frac{\theta_1 - \theta_2}{\sqrt{\theta_1(1-\theta_1) + \theta_2(1-\theta_2)}} \sqrt{n} \right) \hat{f}(\theta_2) d\theta_2 \end{aligned} \quad (35)$$

with $\theta_2 = \theta_1 + \sqrt{2n\theta_1(1-\theta_1)}z$, the expression above then becomes

$$\begin{aligned} & \approx \int_{\frac{-\theta_1\sqrt{n}}{\sqrt{2\theta_1(1-\theta_1)}}}^0 \Phi(z) \hat{f}(\theta_1) \left[\frac{\sqrt{2\theta_1(1-\theta_1)}}{\sqrt{n}} \right] dz \\ & \quad + \int_0^{\frac{(1-\theta_1)\sqrt{n}}{\sqrt{2\theta_1(1-\theta_1)}}} \Phi(-z) \hat{f}(\theta_1) \left[\frac{\sqrt{2\theta_1(1-\theta_1)}}{\sqrt{n}} \right] dz \\ & \approx \underbrace{\left[2 \int_{-\infty}^0 \Phi(z) dz \right]}_{>0} \hat{f}(\theta_1) \frac{\sqrt{2\theta_1(1-\theta_1)}}{\sqrt{n}} \end{aligned} \quad (36)$$

Then we can integrate over θ_1 to obtain:

$$P_{\text{error}} \approx \frac{2}{\sqrt{n\pi}} \int \hat{f}(\theta_1)^2 \sqrt{\theta_1(1-\theta_1)} d\theta_1. \quad (37)$$

This problem can be solved using the method of Lagrange multipliers:

$$\begin{aligned} & \int \sqrt{\theta(1-\theta)} \hat{f}(\theta)^2 d\theta + \lambda \left(\int \hat{f}(\theta) - 1 \right) \\ &= \int (\sqrt{\theta(1-\theta)} \hat{f}(\theta)^2 + \lambda \hat{f}(\theta)) d\theta - \lambda \\ &= \int \mathcal{L}(\theta, \hat{f}, \lambda) d\theta - \lambda \end{aligned} \quad (38)$$

We now calculate the gradient

$$\frac{\partial \mathcal{L}}{\partial \hat{f}} = 2\hat{f}\sqrt{\theta(1-\theta)} + \lambda \quad (39)$$

and then find the optimum for \hat{f} by setting

$$2\hat{f}\sqrt{\theta(1-\theta)} + \lambda = 0 \quad (40)$$

1207

1208

1209

1210

then solving for \hat{f} to obtain

$$\hat{f} = \frac{-\lambda}{2\sqrt{\theta(1-\theta)}}. \quad (41)$$

1211

1212

Taken into consideration our optimization constraint, it can be shown that

1213

1214

$$\int_0^1 \frac{1}{\sqrt{\theta(1-\theta)}} = \frac{1}{\pi}$$

1215

1216

and therefore this implies:

1217

$$\frac{1}{\pi} = \frac{-\lambda}{2}$$

1218

thus requiring:

1219

$$-\lambda = \frac{2}{\pi}.$$

1220

Replacing λ in Eq. 41 we finally obtain

1221

1222

$$\hat{f}(\theta) = \frac{1}{\pi\sqrt{\theta(1-\theta)}} \quad (26 \text{ revisited})$$

1223

1224

Thus the optimal encoding rule is the same (at least in the large- n limit) in this case as when we assume an objective of maximum mutual information (the case considered in **Appendix 1**), though here we assume that the objective is accurate performance of a specific discrimination task.

1225

1226

1227

1228

1229

1230 Appendix 3

1231 Optimal noise for a known prior distribution

1232 Interestingly, we found that the fundamental principles of the theory independently de-
1233 veloped in our work are directly linked to the concept of suprathreshold stochastic reso-
1234 nance (SSR) discovered about two decades ago. Briefly, SSR occurs in an array of n identical
1235 threshold non-linearities, each of which is subject to independently sampled random addi-
1236 tive noise (**Figure 1** in main text). SSR should not be confused with the standard stochastic
1237 resonance (SR) phenomenon. In SR, the amplitude of the input signal is restricted to val-
1238 ues smaller than the threshold for SR to occur. On the other hand, in SSR random draws
1239 from the distribution of input values can exist above threshold levels. Using the simplified
1240 implementational scheme proposed in our work, it can be shown that mutual information
1241 $I(v, k)$ can be also optimized by finding the optimal noise distribution. This is important as
1242 it provides a normative justification as for why sampling must be noisy in capacity-limited
1243 systems. Actually, SSR was initially motivated as a model of neural arrays such as those
1244 synapsing with hair cells in the inner ear, with the direct application of establishing the
1245 mechanisms by which information transmission can be optimized in the design of cochlear
1246 implants (**Stocks et al., 2002**). Our goal in this subsection is to make evident the link be-
1247 tween the novel theoretical implications of our work and the SSR phenomenon developed
1248 in previous work (**Stocks et al., 2002; McDonnell et al., 2007**), which should further justify
1249 our argument of efficient noisy sampling as a general framework for decision behavior, cru-
1250 cially, with a parsimonious implementational nature.

Following our notation, each threshold device (we will call it from now on a *neuron*) can be seen as the number of n resources available to encode an input stimulus v . Here, we assume that each neuron produces a "high" reading if and only if $v + \eta > \tau$, where η is i.i.d. random additive noise (independent of v) following a distribution function f_η , and τ is the minimum threshold required to produce a "high" reading. If we define the noise CDF as F_η , then the probability θ of the neuron giving a "high" reading in response to the input signal v is given by

$$\theta(v) = 1 - F_\eta(\tau - v). \quad (42)$$

It can be shown that the mutual information between the input v and the number of "high" readings k for large n is given by (**McDonnell et al., 2007**)

$$I(v, k) \approx \frac{1}{2} \log_2 \left(\frac{n\pi}{2e} \right) - D_{\text{KL}}[f(v) \| f_J(v)], \quad (43)$$

where f_J is the Jeffreys prior (Eq. 26). Therefore, Jeffreys' prior can also be derived making it a function of the noise distribution f_η

$$f_J(v) = \frac{f_\eta(\tau - v)}{\pi \sqrt{F_\eta(\tau - v)[1 - F_\eta(\tau - v)]}}. \quad (44)$$

Given that the first term in Eq. 43 is always non-negative, a sufficient condition for achieving channel capacity is given by

$$f(v) = f_J(v) \quad \forall v. \quad (45)$$

Typically, the nervous system of any organism has little influence on the distribution of physical signals in the environment. However, it has the ability to shape its internal signals to optimize information transfer. Therefore, a parsimonious solution that the nervous system may adopt to adapt to statistical regularities of environmental signals in a given context is to find the optimal noise distribution f_η^* to achieve channel capacity. Note that this is different from classical problems in communication theory where the goal is usually to find the signal

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

distribution that maximizes mutual information for a channel. Solving Eq. 44 to find $f_{\eta}(v)$ one can find such optimal noise distribution

$$f_{\eta}^*(v) = \frac{\pi}{2} \sin[\pi(1 - F(\tau - v))] f(\tau - v). \quad (46)$$

A further interesting consequence of this set of results is that the ratio between the signal PDF $f(v)$ and the noise PDF f_{η} is

$$\frac{f(v)}{f_{\eta}(\tau - v)} = \frac{2}{\pi \sin[\pi(1 - F(v))]} \quad (47)$$

Using the definition given in Eq. 42 to make this expression a function of θ , one finds the optimal PDF of the encoder

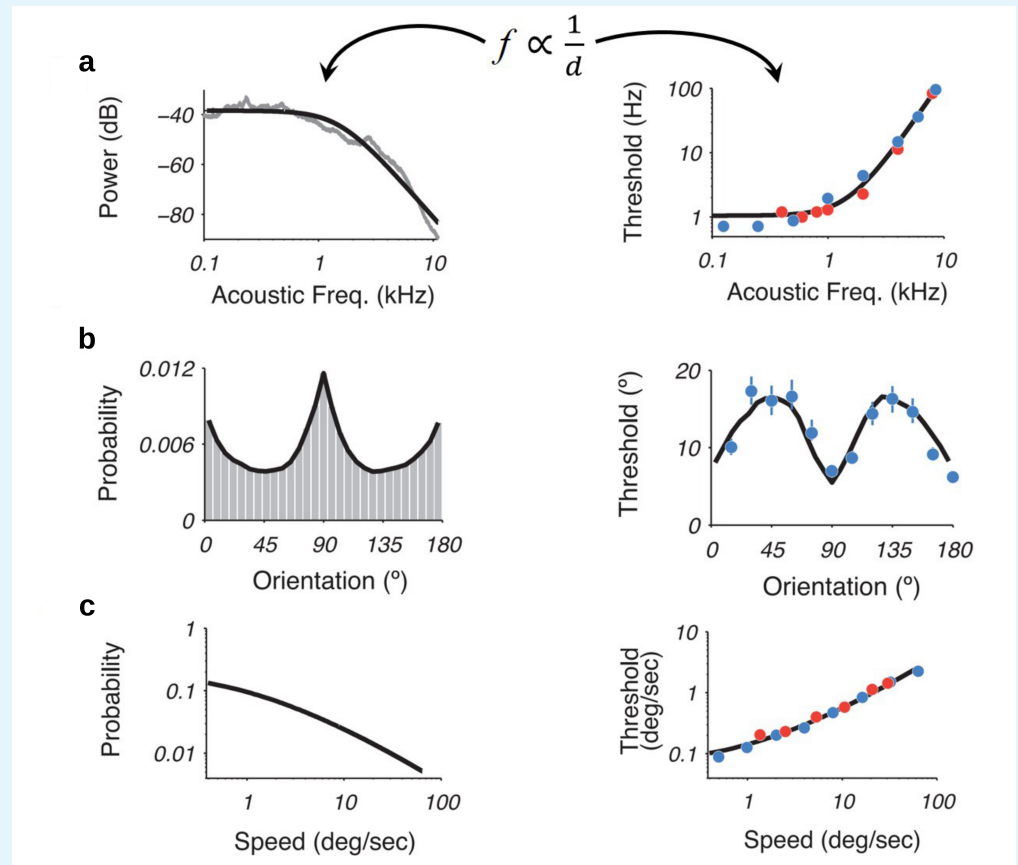
$$f^*(\theta) = \frac{1}{\pi \sqrt{\theta(1 - \theta)}}, \quad (48)$$

which is once again the arcsine distribution (See equations 2 and 5 in main text).

1298 **Appendix 4**

1299
1300

Efficient coding and the relation between environmental priors and discrimination



1301
1302

Appendix 4 Figure 1

1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318

Recently, it was shown that using an efficiency principle for encoding sensory variables, based on population of noisy neurons, it was possible to obtain an explicit relationship between the statistical properties of the environment (the prior) and perceptual discriminability (*Ganguli and Simoncelli, 2016*). The theoretical relation states that discriminability should be inversely proportional to the density of the prior distribution. Interestingly, this relationship holds across several sensory modalities such as **a**) acoustic frequency, **b**) local orientation, **c**) speed (figure adapted with permission from the authors *Ganguli and Simoncelli (2016)*). Here, we investigate whether this particular relation also emerges in our efficient sampling framework.

We first show that we obtain a prediction of exactly the same kind from our model of encoding using a binary channel, in the case that (i) we assume that the encoding rule is optimized for a single environmental distribution, as in the theory of *Ganguli and Simoncelli (2014, 2016)*, and (ii) the objective that is maximized is either mutual information (as in the theory of *Ganguli and Simoncelli*) or the probability of an accurate binary comparison (as considered in *Appendix 2*).

Note that the expected value and variance of a binomial random variable are given by

$$E[r|\theta] = \theta \quad \text{Var}[r|\theta] = \frac{\theta(1-\theta)}{n}, \quad (49)$$

1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373

where we let here $r \equiv k/n$. In **Appendix 2**, we show that if the objective is accuracy maximization, an efficient binomial channel requires that

$$\theta(v) = \left[\sin \left(\frac{\pi}{2} F(v) \right) \right]^2.$$

Thus, replacing $\theta(v)$ in Eq. 49 implies the following relations

$$E[r|\theta] = \sin^2(\omega), \quad \text{Var}[r|\theta] = \frac{\sin^2(\omega)\cos^2(\omega)}{n}, \quad (50)$$

where we let here $\omega \equiv \frac{\pi}{2} F(v)$. Discrimination thresholds d in sensory perception are defined as the ratio between the precision of the representation and the rate of change in the perceived stimulus

$$d \equiv \frac{\sqrt{\text{Var}[r|\theta]}}{E[r|\theta]'}. \quad (51)$$

Substituting the expressions for expected value and variance in Eq. 50 results in

$$\begin{aligned} d &= \frac{1}{2\sqrt{n\omega'}} \\ &= \frac{1}{\sqrt{n\pi} f(v)}. \end{aligned} \quad (52)$$

Thus under our theory, this implies

$$d \propto \frac{1}{f(v)}. \quad (53)$$

This is exactly the relationship derived and tested by **Ganguli and Simoncelli (2016)**.

Our model instead predicts a somewhat different relationship if the encoding rule is required to be robust to alternative possible environmental frequency distributions (the case further discussed in **Appendix 6**). In this case, the robustly optimal encoding rule is DbS, which corresponds to $\theta(v) = F(v)$, rather than the relation 53. Substituting this into Eqs. 49 and 51 yields the prediction

$$d = \frac{\sqrt{F(v)(1-F(v))}}{\sqrt{n}} \cdot \frac{1}{f(v)}. \quad (54)$$

instead of Eq. 52.

One interpretation of the experimental support for the relation 53 reviewed by **Ganguli and Simoncelli (2016)** could be that in the case of early sensory processing of the kind with which they are concerned, perceptual processing is optimized for a particular environmental frequency distribution (representing the long-run experience of an organism or even of the species), so that the assumptions used in **Appendix 2** are the empirically relevant ones. Even so, it is arguable that robustness to changing contextual frequency distributions should be important in the case of higher forms of cognition, so that one might expect prediction 54 to be more relevant for these cases; and indeed, our experimental results for the case of numerosity discrimination are more consistent with Eq. 54 than with 52.

One should also note that even in a case where Eq. 54 holds, if one measures discrimination thresholds over a subset of the stimulus space, over which there is non-trivial variation in $f(v)$, but $F(v)$ does not change very much (because the prior distribution for which the encoding rule is optimized assigns a great deal of probability to magnitudes both higher and lower than those in the experimental data set), then relation (54) restricted to this subset of the possible values for v will imply that the relation (53) should approximately hold. This provides another possible interpretation of the fact that the relation (53) holds fairly well in the data considered by **Ganguli and Simoncelli (2016)**.

1374 **Appendix 5**

1375 **Maximizing expected size of the selected item (fitness maximization)**

1376 We now consider the optimal encoding rule under a different assumed objective, namely,
 1377 maximizing the expected magnitude of the item selected by the subject's response (that is,
 1378 the stimulus judged to be larger by the subject), rather than maximizing the probability of
 1379 a correct response as in **Appendix 2**. While in many perceptual experiments, maximizing
 1380 the probability of a correct response would correspond to maximization of the subject's
 1381 expected reward (or at least maximization of a psychological reward to the subject, who
 1382 is given feedback about the correctness of responses but not about true magnitudes), in
 1383 many of the ecologically relevant cases in which accurate discrimination of numerosity is
 1384 useful to an organism (**Butterworth et al., 2018; Nieder, 2020**), the decision maker's reward
 1385 depends on how much larger one number is than another, and not simply their ordinal rank-
 1386 ing. This would also be true of typical cases in which internal representations of numerical
 1387 magnitudes must be used in economic decision making: the reward from choosing an invest-
 1388 ment with a larger monetary payoff is proportional to the size of the payoff afforded by
 1389 the option that is chosen. Hence it is of interest to consider the optimal encoding rule if we
 1390 suppose that encoding is optimized to maximize performance in a decision task with this
 1391 kind of reward structure.

1392 As in **Appendix 1** and **Appendix 2**, we again consider the problem of optimizing the en-
 1393 coding rule for a specific prior distribution $f(v)$ for the magnitudes that may be encountered,
 1394 and we assume that it is only possible to encode information via "high" or "low" readings.
 1395 The optimization problem that we need to solve is to find the optimal encoding function
 1396 $\theta(v)$ that guarantees a maximal expected value of the chosen outcome, for any given prior
 1397 distribution $f(v)$. Thus the quantity that we seek to maximize is given by

1398
 1399
$$E[v(\text{chosen})] = \int \int f(v_1, v_2) [P_1(\theta(v_1), \theta(v_2))v_1 + P_2(\theta(v_1), \theta(v_2))v_2] dv_1 dv_2 \quad (55)$$

1400 where $P_i(\theta_1, \theta_2)$ is the probability of choosing option i when the encoded values of the two
 1401 options are θ_1 and θ_2 respectively.

1402 We begin by noting that for any pair of input values v_1, v_2 , the integrand in (55) can be
 written as

1403
 1404
$$\begin{aligned} & P_1(\theta(v_1), \theta(v_2))v_1 + P_2(\theta(v_1), \theta(v_2))v_2 \quad (56) \\ &= \max(v_1, v_2) - P_1(\theta(v_1), \theta(v_2))\max(v_2 - v_1, 0) - P_2(\theta(v_1), \theta(v_2))\max(v_1 - v_2, 0) \\ &= \max(v_1, v_2) - [P_1(\theta(v_1), \theta(v_2))I(v_2 > v_1) + P_2(\theta(v_1), \theta(v_2))I(v_1 > v_2)] |v_1 - v_2| \\ &= \max(v_1, v_2) - [P(\text{error} | \theta(v_1), \theta(v_2))I(v_2 > v_1) + P(\text{error} | \theta(v_1), \theta(v_2))I(v_1 > v_2)] |v_1 - v_2| \\ &= \max(v_1, v_2) - P(\text{error} | \theta(v_1), \theta(v_2)) |v_1 - v_2|, \end{aligned}$$

1405 where $I(A)$ is the indicator function (taking the value 1 if statement A is true, and the value
 1406 0 otherwise), and $P(\text{error} | \theta_1, \theta_2)$ is the probability of choosing the lower-valued of the two
 1407 options.

1408 Substituting this last expression for the integrand in (55), we see that we can equivalently
 1409 write

1410
$$E[v(\text{chosen})] = E[\max(v_1, v_2)] - \int \int f(v_1, v_2) P(\text{error} | \theta(v_1), \theta(v_2)) |v_1 - v_2| dv_1 dv_2, \quad (57)$$

where

$$E[\max(v_1, v_2)] \equiv \int \int f(v_1, v_2) \max(v_1, v_2) dv_1 dv_2 \quad (58)$$

1418

1419

1420

1421

1422

1423

1424

is a quantity which is independent of the encoding function $\theta(v)$. Hence choosing $\theta(v)$ to maximize (55) is equivalent to choosing it to minimize

$$E[\text{loss}] = \int \int f(v_1, v_2) P(\text{error} | \theta(v_1), \theta(v_2)) |v_1 - v_2| dv_1 dv_2. \quad (59)$$

As previously specified, the probability of error given two internal noisy readings k_1 and k_2 is given by

$$P(\text{error}) = \left(\frac{k_1}{n} - \frac{k_2}{n} > 0 | v_1, v_2 \right) \quad (60)$$

1425

1426

1427

1428

1429

$$\approx \Phi \left(\frac{\theta_1 - \theta_2}{\sqrt{\frac{\theta_1(1-\theta_1) + \theta_2(1-\theta_2)}{n}}} \right), \quad (61)$$

1430

1431

1432

1433

1434

where in this case we assume that v_1 is the lower-valued option and v_2 is the higher-valued option on any given trial. This implies that $P(\text{error})$ is very close to zero, except when $|\theta_1 - \theta_2| = \mathcal{O}(1/\sqrt{n})$. In this case we have

$$P(\text{error}) \approx \Phi \left(\sqrt{\frac{n}{2}} \frac{\theta_1 - \theta_2}{\sqrt{\theta(1-\theta)}} \right) \quad \text{where } \theta \equiv \frac{\theta_1 + \theta_2}{2}. \quad (62)$$

1435

1436

1437

1438

1439

As in the case of accuracy maximization, here we assume that (v_1, v_2) are independent draws from the same distribution of possible values $f(v)$. Thus $f(v_1, v_2) = f(v_1)f(v_2)$. Then fixing v_1 and integrating over all possible values of v_2 in Eq. 59, the expected loss is approximately

$$E[\text{loss} | v_1] = \int f(v_2) P(\text{error} | v_2, v_1) |v_2 - v_1| dv_2 \quad (63)$$

$$\approx \int f(v_2) \Phi \left(-\sqrt{\frac{n}{2}} \frac{|\theta_1 - \theta_2|}{\sqrt{\theta_1(1-\theta_1)}} \right) |v_2 - v_1| dv_2 \quad (64)$$

$$\approx f(v_1) \int \Phi \left(-\sqrt{\frac{n}{2}} \frac{\theta'(v_1) |v_2 - v_1|}{\sqrt{\theta_1(1-\theta_1)}} \right) |v_2 - v_1| dv_2 \quad (65)$$

$$\approx f(v_1) \int_{-\infty}^{\infty} \Phi(-|z|) \left[\sqrt{\frac{2}{n}} \frac{\theta_1(1-\theta_1)}{\theta'(v_1)} |z| \right] \left[\sqrt{\frac{2}{n}} \frac{\theta_1(1-\theta_1)}{\theta'(v_1)} \right] dz \quad (66)$$

$$\approx \frac{4}{n} \frac{f(v_1)}{\theta'(v_1)^2} [\theta_1(1-\theta_1)] \underbrace{\int_0^{\infty} \Phi(-z) z dz}_{1/4} \quad (67)$$

1440

1441

1442

1443

1444

$$\approx \frac{1}{n} \frac{f(v_1)}{\theta'(v_1)^2} [\theta_1(1-\theta_1)] \quad (68)$$

where in Eq. 66 we have applied the change of variable

$$z \equiv \frac{n}{2} \frac{\theta'(v_1)}{\theta_1(1-\theta_1)} (v_2 - v_1) \quad (69)$$

and in the integral of Eq. 67 we have used

$$\int_0^{\infty} \Phi(-z) z dz = \frac{1}{2} [(z^2 - 1)\Phi(-z) - z\phi(-z)]_0^{\infty} \quad (70)$$

$$= \frac{1}{2} \left[0 - \left(-\frac{1}{2}\right) \right] \quad (71)$$

$$= \frac{1}{4} \quad (72)$$

1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491

where $\phi()$ is the standard normal PDF. Then integrating over v_1 , we have:

$$E[\text{loss}] = \frac{1}{n} \int \frac{f(v_1)^2}{\theta'(v_1)^2} [\theta_1(1 - \theta_1)] dv_1. \quad (73)$$

Thus we want to find the encoding rule $\theta(v)$ to minimize this integral given the prior $f(v)$. We now apply the change of variable $\theta(v) \equiv \sin^2(\gamma(v))$, where $\gamma(v)$ is an increasing function with a range $0 \leq \gamma(v) \leq \frac{\pi}{2}$ for all v . Then we have

$$\theta'(v) = 2 \sin(\gamma(v)) \cos(\gamma(v)) \gamma'(v) \quad (74)$$

$$= 2\sqrt{\theta(v)(1 - \theta(v))} \gamma'(v) \quad (75)$$

and therefore we have

$$\frac{\theta(v)(1 - \theta(v))}{\theta'(v)} = \frac{1}{4} \frac{1}{\gamma'(v)}. \quad (76)$$

This allows us to rewrite Eq. 73 as follows

$$E[\text{loss}] = \frac{1}{n} \int \frac{f(v)^2}{\gamma'(v)^2}. \quad (77)$$

Now the problem is to choose the function $\gamma(v)$ to minimize $E[\text{loss}]$ subject to $0 \leq \gamma(v) \leq \frac{\pi}{2}$. Equivalently, we can choose the function $\gamma'(v) > 0$ to minimize $E[\text{loss}]$ subject to $\int \gamma'(v) dv \leq \frac{\pi}{2}$. Defining $\varphi(v) \equiv \gamma'(v)$, the optimization problem to solve is to choose the function $\varphi(v)$ to

$$\min \int \frac{f(v)^2}{\varphi(v)^2} dv \quad \text{s.t.} \quad \int \varphi(v) dv \leq \frac{\pi}{2} \quad (78)$$

Due to FOC, it can be shown that

$$\frac{f(v)^2}{\varphi(v)^3} = \text{same for all } v \Rightarrow \varphi(v) \sim f(v)^{2/3}. \quad (79)$$

Note also that the constraint $\int \varphi(v) \leq \frac{\pi}{2}$ must hold with equality, thus arriving at

$$\gamma(v) = \frac{\pi}{2} \frac{\int_{-\infty}^v f(\tilde{v})^{2/3} d\tilde{v}}{\int_{-\infty}^{\infty} f(\tilde{v})^{2/3} d\tilde{v}}. \quad (80)$$

Therefore, we finally obtain the efficient encoding rule that maximizes the expected magnitude of the selected item

$$\theta(v) = \sin \left[\frac{\pi}{2} \frac{\int_{-\infty}^v f(\tilde{v})^{2/3} d\tilde{v}}{\int_{-\infty}^{\infty} f(\tilde{v})^{2/3} d\tilde{v}} \right]^2 \quad (81)$$

1492 Appendix 6

1493 **Robust optimality of DbS among Encoding rules with $m = 1$**

1494 Here we consider the nature of the optimal encoding function when the cost of increasing
1495 the size of the sample of values from prior experience that are used to adjust the encoding
1496 rule to the contextual distribution of stimulus values is great enough to make it optimal to
1497 base the encoding of a new stimulus magnitude v on a single sampled value \tilde{v} from the con-
1498 textual distribution. (The conditions required for this to be the case are discussed further
1499 in **Appendix 7**)

1500 We assume that for each of the n independent processing units, the probability of a
1501 "high" reading is given by $\theta(v, \tilde{v}_j)$, where \tilde{v}_j is the draw from the contextual distribution by
1502 processor j , and $\theta(v, \tilde{v})$ is the same function for each of the processing units. The $\{\tilde{v}_j\}$ for $j =$
1503 $1, 2, \dots, n$, are independent draws from the contextual distribution $f(v)$. We further assume
1504 that the function $\theta(v, \tilde{v})$ satisfies certain regularity conditions. First, we assume that θ is a
1505 piecewise continuous function. That is, we assume that the $v - \tilde{v}$ plane can be divided into a
1506 countable number of connected regions, with the boundaries between regions defined by
1507 continuous curves; and that the function $\theta(v, \tilde{v})$ is continuous in the interior of any of these
1508 regions, though it may be discontinuous at the boundaries between regions. And second,
1509 we assume that $\theta(v, \tilde{v})$ is necessarily weakly increasing in v and weakly decreasing in \tilde{v} . The
1510 function is otherwise unrestricted.

1511 For any prior distribution $f(v)$ and any encoding function $\theta(v, \tilde{v})$, we can compute the
1512 probability of an erroneous comparison when two stimulus magnitudes v_1, v_2 are indepen-
1513 dently drawn from the distribution $f(v)$, and each of these stimuli is encoded using n ad-
1514 ditional independent draws $\{\tilde{v}_j\}$ from the same distribution. Let this error probability be
1515 denoted $P_n(\theta; f)$. We wish to find an encoding rule (for given n) that will make this error
1516 probability as small as possible; however, the answer to this question will depend on the
1517 prior distribution $f(v)$. Hence we wish to find an encoding rule that is *robustly* optimal, in
1518 the sense that it achieves the minimum possible value for the upper bound
1519

$$1520 \bar{P}_{error}(\theta) \equiv \sup_{f \in \mathcal{F}} P_n(\theta; f)$$

1521
1522 for the probability of an erroneous comparison. Here the class of possible priors \mathcal{F} to con-
1523 sidered is the set of all possible probability distributions (over values of v) that can be char-
1524 acterized by an integrable probability density function $f(v)$. (We exclude from consideration
1525 priors in which there is an atom of probability mass at some single magnitude v , since in
1526 that case there would be a positive probability of a situation in which it is not clear which
1527 response should be considered "correct", so that P_{error} is not well-defined.) Note that the cri-
1528 terion $\bar{P}_{error}(\theta)$ for ranking encoding rules is not without content, since there exist encoding
1529 rules (including DbS) for which the upper bound is less than 1/2 (the error probability in the
1530 case of a completely uninformative internal representation).

1531 Let us consider first the case in which there is some part of the diagonal line along which
1532 $\tilde{v} = v$ which is not a boundary at which the function $\theta(v, \tilde{v})$ is discontinuous. Then we can
1533 choose an open interval (v_{min}, v_{max}) such that all values v, \tilde{v} with the property that both v and
1534 \tilde{v} lie within the interval (v_{min}, v_{max}) are part of a single region on which $\theta(v, \tilde{v})$ is a continuous
1535 function. Then let θ_{min} be the greatest lower bound with the property that $\theta(v, \tilde{v}) \geq \theta_{min}$ for
1536 all v, \tilde{v} lying within the specified interval, and similarly let θ_{max} be the lowest upper bound
1537 such that $\theta(v, \tilde{v}) \leq \theta_{max}$ for all values within the specified interval. Because of the continuity
1538 of $\theta(v, \tilde{v})$ on this region, as the values v_{min}, v_{max} are chosen to be close enough to each other,
1539 the bounds $\theta_{min}, \theta_{max}$ can be made arbitrarily close to one another.

1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586

Now for any probabilities $0 \leq \theta \leq \theta' \leq 1$, let $P_{\min}(\theta, \theta')$ be the quantity defined in Eq. 30, when $\theta_1 = \theta$ and $\theta_2 = \theta'$; that is, for any v_1, v_2 that are not equal to one another, $P_{\min}(\theta, \theta')$ is the probability of an erroneous comparison if the units representing the smaller magnitude each give a "high" reading with probability θ and those representing the larger magnitude each give a "high" reading with probability θ' . Then the probability of erroneous choice P_{error} when $f(v)$ is a distribution with support entirely within the interval (v_{\min}, v_{\max}) is necessarily greater than or equal to the lower bound $P_{\min}(\theta_{\min}, \theta_{\max})$. The reason is that for any v_1, v_2 in the support of $f(v)$, the probabilities

$$\theta_i = \int \theta(v_i, \bar{v}) f(\bar{v}) d\bar{v}$$

will necessarily lie within the bounds $\theta_{\min} \leq \theta_i \leq \theta_{\max}$ for both $i = 1, 2$. Given these bounds, the most favorable case for accurate discrimination between the two magnitudes will be to assign the largest possible probability θ_{\max} to units being on in the representation of the larger magnitude, and the smallest possible probability θ_{\min} to units being on in the representation of the smaller magnitude. Since the lower bound $P_{\min}(\theta_{\min}, \theta_{\max})$ applies in the case of any individual values v_1, v_2 drawn from the support of $f(v)$, this same quantity is also a lower bound for the average error rate integrating over the prior distributions for v_1 and v_2 .

One can also show that as the two bounds $\theta_{\min}, \theta_{\max}$ approach one another, the lower bound $P_{\min}(\theta_{\min}, \theta_{\max})$ approaches 1/2, regardless of the common value that θ_{\min} and θ_{\max} both approach. Hence it is possible to make $P_{\min}(\theta_{\min}, \theta_{\max})$ arbitrarily close to 1/2, by choosing values for v_{\min}, v_{\max} that are close enough to one another. It follows that for any bound P_{\min} less than 1/2 (including values arbitrarily close to 1/2), we can choose a prior distribution $f(v)$ for which P_{error} is necessarily equal to P_{\min} or larger. It follows that in the case of a function $\theta(v, \bar{v})$ of this kind, the upper bound $\bar{P}_{\text{error}}(\theta)$ is equal to 1/2.

In order to achieve an upper bound lower than 1/2, then, we must choose a function $\theta(v, \bar{v})$ that is discontinuous along the entire line $v = \bar{v}$. For any such function, let us consider a value v^* with the property that all points (v, \bar{v}) near (v^*, v^*) with $v > \bar{v}$ belong to one region on which θ is continuous, and all points near (v^*, v^*) with $v < \bar{v}$ belong to another region. Then under the assumption of piecewise continuity, $\theta(v, \bar{v})$ must approach some value $\bar{\theta}(v^*)$ as the values (v, \bar{v}) converge to (v^*, v^*) from within the region where $v > \bar{v}$, and similarly $\theta(v, \bar{v})$ must approach some value $\underline{\theta}(v^*)$ as the values (v, \bar{v}) converge to (v^*, v^*) from within the region where $v < \bar{v}$.

It must also be possible to choose values $v_{\min} < v^* < v_{\max}$ such that all points (v, v) with $v_{\min} < v < v_{\max}$ are points on the boundary between the two regions on which θ is continuous. Given such values, we can then define bounds $\underline{\theta}_{\min}, \underline{\theta}_{\max}, \bar{\theta}_{\min},$ and $\bar{\theta}_{\max}$, such that

$$\underline{\theta}_{\min} \leq \theta(v, \bar{v}) \leq \underline{\theta}_{\max}$$

for all $v_{\min} < v < \bar{v} < v_{\max}$, and

$$\bar{\theta}_{\min} \leq \theta(v, \bar{v}) \leq \bar{\theta}_{\max}$$

for all $v_{\min} < \bar{v} < v < v_{\max}$. Moreover, piecewise continuity of the function $\theta(v, \bar{v})$ implies that by choosing both v_{\min} and v_{\max} close enough to v^* we can make the bounds $\underline{\theta}_{\min}, \underline{\theta}_{\max}$ arbitrarily close to $\underline{\theta}(v^*)$, and make the bounds $\bar{\theta}_{\min}, \bar{\theta}_{\max}$ arbitrarily close to $\bar{\theta}(v^*)$.

Next, for any set of four probabilities $0 \leq \underline{\theta} \leq \underline{\theta}' \leq 1$ and $0 \leq \bar{\theta} \leq \bar{\theta}' \leq 1$, let us define

$$\hat{P}_{\min}(\underline{\theta}, \underline{\theta}'; \bar{\theta}, \bar{\theta}') \equiv \text{E}[P_{\min}(\theta(z_1), \theta'(z_2)) | z_1 < z_2], \quad (82)$$

where

$$\theta(z) \equiv z\bar{\theta} + (1-z)\underline{\theta}, \quad \theta'(z) \equiv z\bar{\theta}' + (1-z)\underline{\theta}', \quad (83)$$

1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649

and z_1, z_2 are two independent random variables, each distributed uniformly on $[0, 1]$. Then if $\theta(v, \bar{v})$ lies between the lower bound $\underline{\theta}$ and upper bound $\underline{\theta}'$ whenever $v < \bar{v}$, and between the lower bound $\bar{\theta}$ and upper bound $\bar{\theta}'$ whenever $v > \bar{v}$, then the probability θ of a processing unit representing the magnitude v giving a "high" reading will lie between the bounds $\theta(z) \leq \theta \leq \theta'(z)$, where $z = F(v)$ is the quantile of v within the prior distribution. It follows that in the case of any two magnitudes v_1, v_2 with $v_1 < v_2$, the probability of an erroneous comparison will be bounded below by $P_{\min}(\theta(z_1), \theta'(z_2))$, where $z_i = F(v_i)$ for $i = 1, 2$, since the probability of a correct discrimination will be maximized by making the units representing v_1 give as few high readings as possible and the units representing v_2 give as many high readings as possible. Integrating over all possible draws of v_1, v_2 , one finds that the quantity $\hat{P}_{\min}(\underline{\theta}, \underline{\theta}'; \bar{\theta}, \bar{\theta}')$ defined in (82) is a lower bound for the overall probability of an erroneous comparison, given that regardless of the prior $f(v)$, the quantiles z_1, z_2 will be two independent draws from the uniform distribution on $[0, 1]$.

Now consider again an encoding function $\theta(v, \bar{v})$ of the kind discussed two paragraphs above, and an interval of stimulus values (v_{\min}, v_{\max}) of the kind discussed there. For any prior distribution $f(v)$ with support entirely contained within the interval (v_{\min}, v_{\max}) , the probability of an erroneous comparison is bounded below by

$$P_n(\theta; f) \geq \hat{P}_{\min}(\underline{\theta}_{\min}, \underline{\theta}_{\max}; \bar{\theta}_{\min}, \bar{\theta}_{\max}),$$

where the function \hat{P}_{\min} is defined in (82). Moreover, by choosing the values v_{\min}, v_{\max} close enough to v^* , we can make this lower bound arbitrarily close to $P^e(\underline{\theta}(v^*), \bar{\theta}(v^*))$, where for any probabilities $\underline{\theta}, \bar{\theta}$ we define

$$P^e(\underline{\theta}, \bar{\theta}) \equiv \hat{P}_{\min}(\underline{\theta}, \underline{\theta}; \bar{\theta}, \bar{\theta}). \quad (84)$$

Hence in the case of the encoding function considered, the upper bound $\bar{P}_{\text{error}}(\theta)$ must be at least as large as $P^e(\underline{\theta}(v^*), \bar{\theta}(v^*))$. We further observe that the quantity $P^e(\underline{\theta}, \bar{\theta})$ defined in (84) is just the probability of an erroneous comparison in the case of an encoding rule according to which

$$\begin{aligned} \theta(v, \bar{v}) &= \underline{\theta} & \text{if } v < \bar{v}, \\ \theta(v, \bar{v}) &= \bar{\theta} & \text{if } v > \bar{v}. \end{aligned}$$

Note that in the case of such an encoding rule, the probability of an erroneous comparison is the same for all prior distributions, since under this rule all that matters is the distribution of the quantile ranks of v and \bar{v} . It is moreover clear that $P^e(\underline{\theta}, \bar{\theta})$ is an increasing function of $\underline{\theta}$ and a decreasing function of $\bar{\theta}$. It thus achieves its minimum possible value if and only if $\underline{\theta} = 0$ and $\bar{\theta} = 1$, in which case it takes the value P_{error}^{DbS} , the probability of erroneous comparison in the case of decision by sampling (again, independent of the prior distribution).

Thus in the case that there exists any magnitude v^* for which $\underline{\theta}(v^*) > 0$, $\bar{\theta}(v^*) < 1$, or both, there exist priors $f(v)$ for which $P_n(\theta; f)$ must exceed $P_{\text{error}}^{DbS} = P^e(0, 1)$. Hence in order to minimize the upper bound $\bar{P}_{\text{error}}(\theta)$, it must be the case that $\underline{\theta}(v) = 0$ and $\bar{\theta}(v) = 1$ for all v . But then our assumption that the encoding rule $\theta(v, \bar{v})$ is at least weakly increasing in v and at least weakly decreasing in \bar{v} requires that

$$\begin{aligned} \theta(v, \bar{v}) &= 0 & \text{for all } v < \bar{v}, \\ \theta(v, \bar{v}) &= 1 & \text{for all } v > \bar{v}. \end{aligned}$$

Thus the encoding rule must be the DbS rule, the unique rule for which $\bar{P}_{\text{error}}(\theta)$ is no greater than P_{error}^{DbS} .

1650 **Appendix 7**

1651 **Sufficient conditions for the optimality of Dbs**

1652 Here we consider the general problem of choosing a value of m (the number of samples
1653 from the contextual distribution $f(v)$ to use in encoding any individual stimulus) and an
1654 encoding rule $\theta(v; \tilde{v}_1, \dots, \tilde{v}_m)$ to be used by each of the n processing units that encode the
1655 magnitude of that single stimulus, so as to minimize the compound objective
1656

$$1657 \quad \bar{P}_{error}(\theta) + K(m),$$

1658 where \bar{P}_{error} is the upper bound on the probability of an erroneous comparison under the
1659 encoding rule θ , and $K(m)$ is the cost of using a sample of size m when encoding each stimu-
1660 lus magnitude. The value of n is taken as fixed at some finite value. (This too can be
1661 optimized subject to some cost of additional processing units, but we omit formal analysis
1662 of this problem.) We assume that $K(m)$ is an increasing function of m , and can without loss
1663 of generality assume the normalization $K(0) = 0$. In this optimization problem, we assume
1664 that the only encoding functions θ to be considered are ones that are piecewise continuous,
1665 at least weakly increasing in v , and weakly decreasing in each of the \tilde{v}_j .

1666 For any value of m , let $P^*(m)$ be the minimum achievable value for $\bar{P}_{error}(\theta)$. (**Appendix 6**
1667 illustrates how this kind of problem can be solved, for the case $m = 1$.) Then the optimal
1668 value of m will be the one that minimizes $P^*(m) + K(m)$.

1669 We can establish a lower bound for $P^*(m)$ that holds for any m :

$$1670 \quad \begin{aligned} P^*(m) &\equiv \inf_{\theta(v; \tilde{v}_1, \dots, \tilde{v}_m)} \sup_{f \in \mathcal{F}} P_n(\theta; f) \\ &\geq \sup_{f \in \mathcal{F}} \inf_{\theta(v; \tilde{v}_1, \dots, \tilde{v}_m)} P_n(\theta; f) \\ &= \sup_{f \in \mathcal{F}} \inf_{\theta(v)} P_n(\theta; f) \equiv \underline{P}_n. \end{aligned} \quad (85)$$

1671 In the second line, we allow the function $\theta(v; \tilde{v}_1, \dots, \tilde{v}_m)$ to be chosen after a particular prior
1672 $f(v)$ has already been selected, which cannot increase the worst-case error probability. In
1673 the third line, we note that the only thing that matters about the encoding function chosen in
1674 the second line is the mean value of $\theta(v; \tilde{v}_1, \dots, \tilde{v}_m)$ for each possible magnitude v , integrating
1675 over the possible samples of size m that may be drawn from the specified prior; hence we
1676 can more simply write the problem on the second line as one involving a direct choice of
1677 a function $\theta(v)$, which may be different depending on the prior $f(v)$ that has been chosen.
1678 The problem on the third line defines a bound \underline{P}_n that does not depend on m .

1679 A set of sufficient conditions for $m = 1$ to be optimal is then given by the assumptions
1680 that

- 1681 (a) $P^*(0) > P^*(1) + K(1)$, and
1682 (b) $P^*(1) - \underline{P} < K(2) - K(1)$.

1683 Condition (a) implies that $m = 0$ will be inferior to $m = 1$: the cost of a single sample is not so
1684 large as to outweigh the reduction in $\bar{P}_{error}(\theta)$ that can be achieved using even one sample.
1685 Condition (b) implies that $m = 1$ will be superior to any $m' > 1$. The lower bound (85), together
1686 with our monotonicity assumption regarding $K(m)$, implies that for any $m' > 1$,

$$1687 \quad P^*(1) - P^*(m') \leq P^*(1) - \underline{P} < K(2) - K(1) \leq K(m') - K(1),$$

1688 and hence that

$$1689 \quad P^*(1) + K(1) < P^*(m') + K(m').$$

1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744

While condition (b) is stronger than is needed for this conclusion, the sufficient conditions stated in the previous paragraph have the advantage that we need only consider optimal encoding rules for the cases $m = 0$ and $m = 1$, and the efficient coding problem stated in definition (85), in order to verify that the conditions are both satisfied. The efficient coding problem for the case $m = 1$ is treated in **Appendix 6**, where we show that $P^*(1) = P_{error}^{DbS} < 1/2$. Using the calculations explained in **Appendix 2**, we can provide an analytical approximation to this quantity in the limiting case of large n .

Equation 37 states that for any encoding rule $\theta(v)$ and any prior distribution $f(v)$, the value of P_{error} for any large enough value of n will approximately equal

$$P_n(\theta; f) \approx \frac{2}{\sqrt{n\pi}} \int \hat{f}(\tilde{\theta})^2 \sqrt{\tilde{\theta}(1-\tilde{\theta})} d\tilde{\theta}, \quad (37 \text{ revisited})$$

where $\hat{f}(\theta)$ is the probability density function of the distribution of values for $\theta(v)$ implied by the function $\theta(v)$ and the distribution $f(v)$ of values for v . In the case of DbS, the probability distribution over alternative internal representations k_i (and hence the probability of error) is the same as in the case of an encoding rule $\theta(v) = F(v)$, so that equation 37 can be applied. Furthermore, for any prior distribution $f(v)$, the probability distribution of values for the quantile $z = F(v)$ will be a uniform distribution over the interval $[0, 1]$, so that $\hat{f}(\theta) = 1$ for all θ . It follows that

$$P_{error}^{DbS, \text{lim}} \approx \frac{2}{\sqrt{n\pi}} \int \sqrt{\tilde{\theta}(1-\tilde{\theta})} d\tilde{\theta} = \frac{1}{4} \sqrt{\frac{\pi}{n}}. \quad (86)$$

In the case that $m = 0$, instead, the same function $\theta(v)$ must be used regardless of the contextual distribution $f(v)$. Under the assumption that $\theta(v)$ is piecewise continuous, there must exist a magnitude v^* such that $\theta(v)$ is continuous over some interval (v_{min}, v_{max}) containing v^* in its interior. Let $\theta_{min}, \theta_{max}$ be the greatest lower bound and least upper bound respectively, such that

$$\theta_{min} \leq \theta(v) \leq \theta_{max}$$

for all $v_{min} < v < v_{max}$. The continuity of $\theta(v)$ on this interval means that by choosing both v_{min} and v_{max} close enough to v^* , we can make both θ_{min} and θ_{max} arbitrarily close to $\theta(v^*)$.

By the same argument as in **Appendix 6**, for any prior distribution $f(v)$ with support entirely contained in the interval (v_{min}, v_{max}) , the pair of stimulus magnitudes v_1, v_2 will have to imply $\theta_{min} \leq \theta(v_1), \theta(v_2) \leq \theta_{max}$ with probability 1, and as a consequence the error probability $P_n(\theta; f)$ will necessarily be greater than or equal to the lower bound $P_{min}(\theta_{min}, \theta_{max})$. By choosing both v_{min} and v_{max} close enough to v^* , we can make this lower bound arbitrarily close to $P_{min}(\theta(v^*), \theta(v^*)) = 1/2$. Hence for any encoding rule $\theta(v)$ with $m = 0$, the upper bound $\bar{P}_{error}(\theta)$ cannot be lower than 1/2. It follows that $P^*(0) = 1/2$.

Given this, condition (a) can alternatively be expressed as

$$P_{error}^{DbS} + K(1) < 1/2.$$

Note that if $K(1)$ remains less than 1/2 no matter how large n is, this condition will necessarily be satisfied for all large enough values of n , since (86) implies that P_{error}^{DbS} eventually becomes arbitrarily small, in the case of large enough n . (On the other hand, the condition can easily be satisfied for some range of smaller values of n , even if $K(1) > 1/2$ once n becomes very large.)

In order to consider the conditions under which condition (b) will also be satisfied, it is necessary to further analyze the efficient coding problem stated in (85). We first observe that for any prior $f(v) \in \mathcal{F}$ and encoding rule $\theta(v)$, the encoding rule can always be expressed in the form $\theta(v) = \varphi(F(v))$, where $\varphi(z)$ is a piecewise-continuous, weakly increasing function giving the probability of a "high" reading as a function of the quantile z of the stimulus

1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1763
1762
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780

1781
1782
1783
1784
1785
1786
1787
1788

magnitude in the prior distribution. We then note that when this representation is used for the encoding function in problem 85, the error probability $P_n(\theta; f)$ depends only on the function $\varphi(z)$, in a way that is independent of the prior $f(v)$. Hence the inner minimization problem in Eq. 85 can equivalently be written as

$$\inf_{\varphi(z)} P_n(\varphi). \quad (87)$$

This problem has a solution for the optimal $\varphi(z)$ for any number of processing units n , and an associated value, that is independent of the prior $f(v)$. Hence we can write the bound defined in (85) more simply as

$$\underline{P}_n = \inf_{\varphi(z)} P_n(\varphi). \quad (88)$$

Condition (b) will be satisfied as long as the bound defined in (88) is not too much lower than $P_{\text{error}}^{\text{DbS}}$. In fact, this bound can be a relatively large fraction of $P_{\text{error}}^{\text{DbS}}$. We consider the problem of the optimal choice of an encoding function $\theta(v)$ for a known prior $f(v)$ in **Appendix 2**. In the limiting case of a sufficiently large n , substitution of equation 2 into 37 yields the approximate solution

$$\underline{P}_n^{\text{lim}} \approx \frac{2}{\sqrt{n\pi}} \frac{1}{\pi^2} \frac{d\tilde{\theta}}{\sqrt{\tilde{\theta}(1-\tilde{\theta})}} = \frac{2}{\sqrt{n\pi^3}}. \quad (89)$$

Thus as n is made large, the ratio $\underline{P}_n^{\text{lim}} / P_{\text{error}}^{\text{DbS,lim}}$ converges to the value

$$\underline{P}_n^{\text{lim}} / P_{\text{error}}^{\text{DbS,lim}} = 8/\pi^2 = 0.81. \quad (90)$$

This means that increases in the sample size m above 1 cannot reduce $P^*(m)$ by even 20 percent relative to $P^*(1)$, no matter how large the sample may be, whereas $P^*(1)$ may be only a small fraction of $P^*(0)$ (as is necessarily the case when n is large). This makes it quite possible for $K(2) - K(1)$ to be larger than $P_{\text{error}}^{\text{DbS}} - \underline{P}$ while at the same time $P^*(0) - P_{\text{error}}^{\text{DbS}}$ is larger than $K(1)$. **In this case, the optimal sample size will be $m = 1$, and the optimal encoding rule will be DbS.**

While these analytical results for the asymptotic (large- n) case are useful, we can also numerically estimate the size of the terms $P^*(0)$, \underline{P} , and $P_{\text{error}}^{\text{DbS}}$ in the case of any finite value for n . We have derived an exact analytical value for $P^*(0) = 1/2$ above. The quantity $P_{\text{error}}^{\text{DbS}}$ can be computed through Monte Carlo simulation for any value of n . (Note that this calculation depends only on n , and is independent of the contextual distribution $f(v)$; we need only to calculate $P_n(\varphi)$ for the function $\varphi(z) = z$.) The calculation of \underline{P}_n for a given finite value of n is instead more complex, since it requires us to optimize $P_n(\varphi)$ over the entire class of possible functions $\varphi(z)$.

Our approach is to estimate the minimum achievable value of $P_n(\varphi)$ by finding the minimum achievable value over a flexible parametric family of possible functions $\varphi(z)$. We specify the function φ in terms of the implied $\hat{F}(\theta)$, the CDF for values of $\theta(v)$. We let $\hat{F}(\theta)$ be implicitly defined by

$$[\sin((\pi/2)\hat{F}(\theta))]^2 = g(\theta), \quad (91)$$

where $g(\theta)$ is a function of θ with the properties that $g(0) = 0$, $g(1) = 1$, as required for $\hat{F}(\theta)$ to be the CDF of a probability distribution. More specifically, we assume that $g(\theta)$ is a finite-order polynomial function consistent with these properties, which require that it can be written in the form

$$g(\theta) = \theta [1 + (\theta - 1) (g_0 + g_1\theta + \dots + g_p\theta^p)], \quad (92)$$

where $\{g_0, \dots, g_p\}$ are a set of parameters over which we optimize. Note that for a large enough value of p , any smooth function can be well approximated by a member of this family. At the same time, our choice of a parametric family of functions has the virtue that

1803

1805

1806

1807

1808

1809

1810

1811

1812

1813

1814

1815

1816

1817

1818

1819

1820

1821

1822

1823

1824

1825

1826

1827

1828

1830

1831

1832

1833

1834

1835

1836

1837

1838

1839

1840

1841

1842

1843

1844

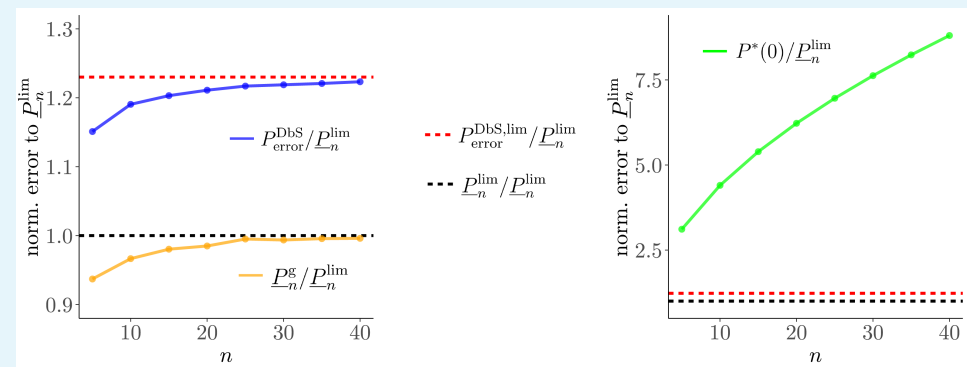
the CDF that corresponds to the optimal coding rule in the large- n limit belongs to this family (regardless of the value of p), since this coding rule (equation 3) corresponds to the case $g_0 = \dots = g_p = 0$ of equation 92.

We computed via numerical simulations the best encoder function assuming $g(\theta)$ to be of order 5 (Eq. 92) for various finite values of $n = [5, 10, 15, 20, 25, 30, 35, 40]$, and we define the expected error of this optimal encoder for a given n to be \underline{P}_n^g (i.e., a lower bound for P_n within the family of functions defined by g). Our goal is to compare this quantity to the asymptotic approximation $\underline{P}_n^{\text{lim}}$, in order to evaluate how accurate the asymptotic approximation is.

Additionally, we also compute the value $P_{\text{error}}^{\text{DbS}}$ for each finite value of n through Monte Carlo simulation (please note that $P_{\text{error}}^{\text{DbS}}$ is different from the quantity $P_{\text{error}}^{\text{DbS,lim}}$ defined in Eq. 86, that is only an asymptotic approximation for large n). Then, we can compare $P_{\text{error}}^{\text{DbS}}$ to the value predicted by the asymptotic approximations $P_{\text{error}}^{\text{DbS,lim}}$ and $\underline{P}_n^{\text{lim}}$.

Another quantity that is important to compute, in order to determine whether DbS can be optimal when n is not too large, is the size of $P^*(0)$ relative to the quantities computed above. Since $P^*(0)$ does not shrink as n increases, it is obvious that $P^*(0)$ is much larger than the other quantities in the large- n limit. But how much bigger is it when n is small? To investigate this, we compute the value of the ratio $P^*(0)/\underline{P}_n^{\text{lim}}$ when n is small. This quantity is given by

$$\frac{P^*(0)}{\underline{P}_n^{\text{lim}}} = \frac{\sqrt{4n\pi^3}}{4} \quad (93)$$



Appendix 7 Figure 1

In Appendix 7-Figure 1, all error quantities discussed above are normalized relative to $\underline{P}_n^{\text{lim}}$. The black dashed lines in both panels represent $(\underline{P}_n^{\text{lim}}/\underline{P}_n^{\text{lim}}) = 1$. The ratio of the asymptotic approximation for $P_{\text{error}}^{\text{DbS,lim}}$ relative to $\underline{P}_n^{\text{lim}}$ is plotted with the red dashed lines, where $(P_{\text{error}}^{\text{DbS,lim}}/\underline{P}_n^{\text{lim}}) \approx 1.23$. Note that the sufficient conditions for DbS to be optimal can be stated as

- (a) $K(1) < P^*(0) - P_{\text{error}}^{\text{DbS}}$, and
- (b) $K(2) - K(1) > P_{\text{error}}^{\text{DbS}} - \underline{P}_n$.

Therefore, Appendix 7-Figure 1 shows the numerical magnitudes of the expressions on the right-hand side of both inequalities (normalized by the value of $\underline{P}_n^{\text{lim}}$). The most important result from the analyses presented in this figure is that even for small values of n , the right-hand side of the first inequality (see right panel) will be a much larger quantity than the right-hand side of the second inequality (see left panel). Thus it can easily be the case that $K(1)$ and $K(2)$ are such that both inequalities are satisfied: it is worth increasing m from 0 to 1, but not worth increasing m to any value higher than 1. In this case, the optimal sample size will be $m = 1$, and the optimal encoding rule will be DbS.

1845 Additionally, we found that the computations of $P_{\text{error}}^{\text{DbS}}$ for each finite value of n are slightly
1846 higher than $\underline{P}_n^{\text{lim}}$ even for small n values (blue line in the left panel), but quickly reach the
1847 asymptotic value $P_{\text{error}}^{\text{DbS,lim}} / \underline{P}_n^{\text{lim}}$ as n increases. Thus, even for small values of n , the asymp-
1848 totic approximation of optimal performance for the case of complete prior knowledge is
1849 superior than DbS. We also found that the computations of \underline{P}_n^g for each finite value of n can-
1850 not reduce $\underline{P}_n^{\text{lim}}$ by even 5 percent for small n values (orange line in the left panel). Moreover,
1851 \underline{P}_n^g quickly reached the asymptotic value $\underline{P}_n^{\text{lim}}$, thus suggesting that the asymptotic solution
1852 is virtually indistinguishable from the optimal solution (at least based on the flexible fam-
1853 ily of g functions) also for finite values of n , which crucially are in the range of the values
1854 found to explain the data in the numerosity discrimination experiment of our study. Thus,
1855 these results confirm that the asymptotic approximations used in our study are not likely
1856 to influence the conclusions of the experimental data in our work.

1857 Appendix 8

1858 **Relation to Bhui and Gershman (2018)**

1859 **Bhui and Gershman (2018)** also argue that an efficient coding scheme can be implemented
1860 by a version of DbS. However, both the efficient coding problem that they consider, and the
1861 version of DbS that they consider, are different than in our analysis, so that our results are
1862 not implied by theirs.

1863 Like us, Bhui and Gershman consider encoding schemes in which the internal representa-
1864 tion r must take one of a finite number of values. However, their efficient coding problem
1865 considers the class of all encoding rules that assign one or another of N possible values
1866 of r to a given stimulus v . In their discussion of the ideal efficient coding benchmark, they
1867 do not require r to be the ensemble of output states of a set of n neurons, each of which
1868 must use the same rule as the other units, and therefore consider a more flexible family of
1869 possible encoding rules, as we explain in more detail below.

1870 The encoding rule that solves our efficient coding problem is stochastic; even under the
1871 assumption that the prior $f(v)$ is known with perfect precision (the case of unbounded m in
1872 the more general specification of our framework, so that sampling error in estimation of this
1873 distribution from prior experience is not an issue), we show that it is optimal for the prob-
1874 abilities $p(k|v)$ not to all equal either zero or one. The optimal rule within the more flexible
1875 class considered by Bhui and Gershman is instead deterministic: each stimulus magnitude
1876 v is assigned to exactly one category k with certainty. The boundaries between the set of
1877 $n + 1$ categories furthermore correspond to the quantiles $(1/(n + 1), 2/(n + 1), \dots, n/(n + 1))$ of
1878 the prior distribution, so that each category is used with equal frequency. Thus the optimal
1879 encoding rule is given by a deterministic function $y(v)$, a non-decreasing step function that
1880 takes $n + 1$ discrete values.

1881 Bhui and Gershman show that when there is no bound on m , the number of samples
1882 from prior experience that can be used to estimate the contextual distribution — their opti-
1883 mal encoding rule for a given number of categories N — can be implemented by a form
1884 of DbS. However, the DbS algorithm that they describe is different than in our discussion.
1885 Bhui and Gershman propose to implement the deterministic classification $y(v)$ by comput-
1886 ing the fraction of the sampled values \tilde{v} that are less than v . In the limiting case of an infinite
1887 sample from the prior distribution, this fraction is equal to $F(v)$ with probability one, and
1888 $y(v)$ is then determined by which of the intervals $[0, 1/N), [1/N, 2/N), \dots, [(N - 1)/N, 1]$ the
1889 quantile $F(v)$ falls within. Thus whereas in our discussion, DbS is an algorithm that allows
1890 each of our units to compute its state using only a single sampled value \tilde{v}_j , the DbS algo-
1891 rithm proposed by Bhui and Gershman to implement efficient coding is one in which a large
1892 number of sampled values are used to jointly compute the output states of all of the units
1893 in a coordinated way.

Bhui and Gershman also consider the case in which only a finite number of samples
($\tilde{v}_1, \dots, \tilde{v}_m$) can be used to compute the representation k_i of a given stimulus magnitude v_i ,
and ask what kind of rule is efficient in that case. They show that in this case a variant of
DbS with kernel-smoothing is superior to the version based on the empirical quantile of v_i
(which now involves sampling error). In this more general case, the variant DbS algorithms
considered by Bhui and Gershman make the representation k_i of a given stimulus proba-
bilistic; but the class of probabilistic algorithms that they consider remains different from
the one that we discuss. In particular, they continue to consider algorithms in which the
category k_i can be an arbitrary function of v_i and a single set of m sampled values that is
used to compute the complete representation; they do not impose the restriction that k_i
be the number of units giving a "high" reading when the output state of each of n individ-

1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921

ual processing units is computed independently using the same rule (but an independent sample of values from prior experience in the case of each unit).

The kernel-smoothing algorithms that they consider are based on a finite set of m pairwise comparisons between the stimulus magnitude v_i and particular sampled values \tilde{v}_j , the outcomes of which are then aggregated to obtain the internal representation k_i . However, they allow the quantity $K(v_i - \tilde{v}_j)$ computed by comparing v_i to an individual sampled value to vary continuously between 0 and 1, rather than having to equal either 0 or 1, as in our case (where the state of an individual unit must be either "high" or "low"). The quantities $K(v_i - \tilde{v}_j)$ are able to be summed with perfect precision, before the resulting sum is then discretized to produce a final representation that takes one of only N possible values. Thus an assumption that only finite-precision calculations are possible is made only at the stage where the final output of the joint computation of the processors must be "read out"; the results of the individual binary comparisons are assumed to be integrated with infinite precision. In this respect, the algorithms considered by Bhui and Gershman are not required to economize on processing resources in the same sense as the class that we consider; the efficient coding problem for which they present results is correspondingly different from the problem that we discuss for the case in which m is finite.

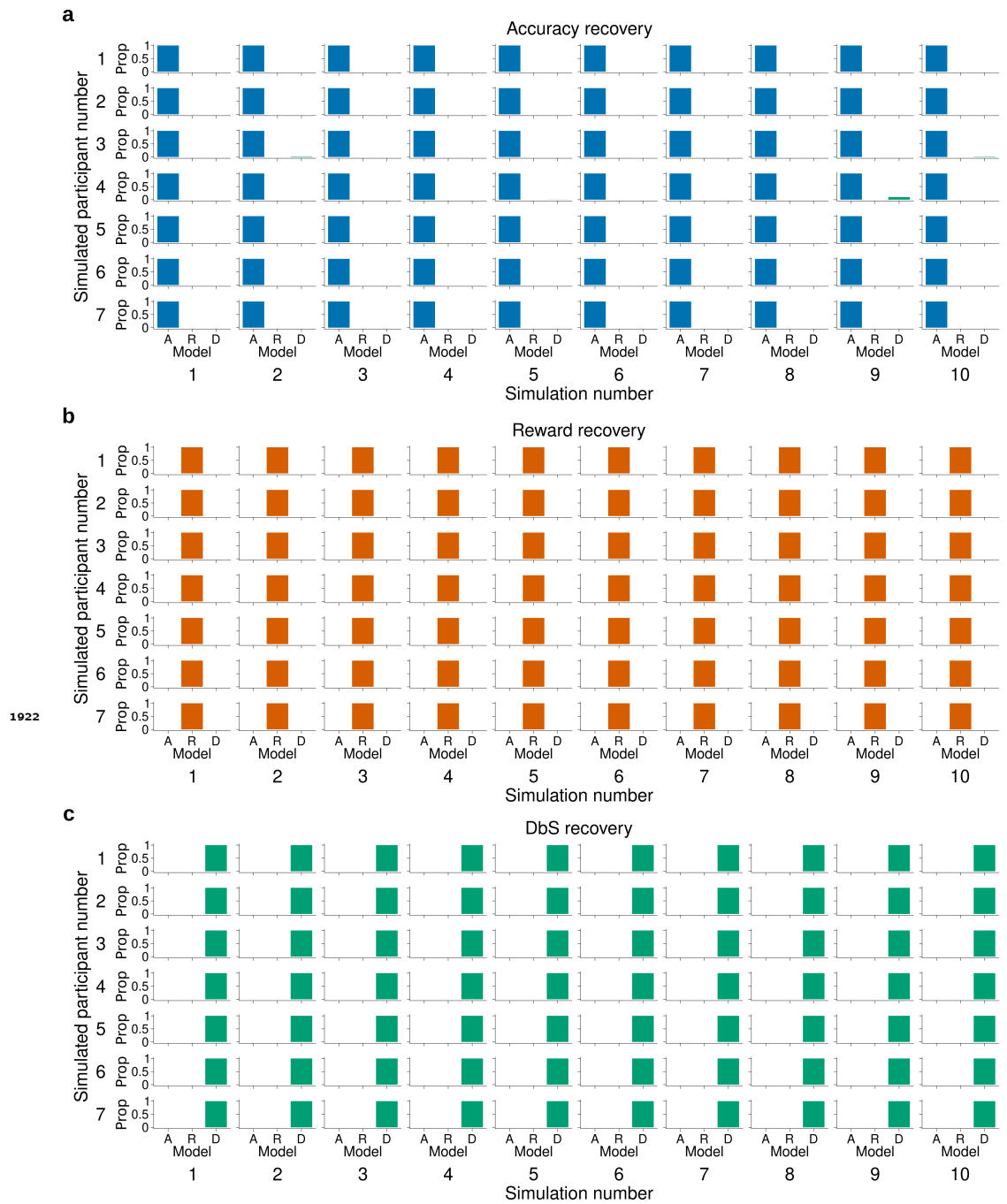


Figure 3-Figure supplement 1. Model recovery for α fixed. The latent-mixture model was fit to synthetic data obtained by simulating 10 times each encoding rule on the trials from participants of Experiment 1. This also means that we used the same number of trials per condition that each participant experienced in our experiments. Each histogram shows the proportion (Prop) of the recovered encoding rule for synthetic data from **a)** the accuracy maximizing encoding rule θ_A , **b)** the reward maximizing encoding rule θ_R , and **c)** decision by sampling θ_D . The latent mixture model can accurately recover the underlying encoding rule. In this model the α parameter was set to 2.

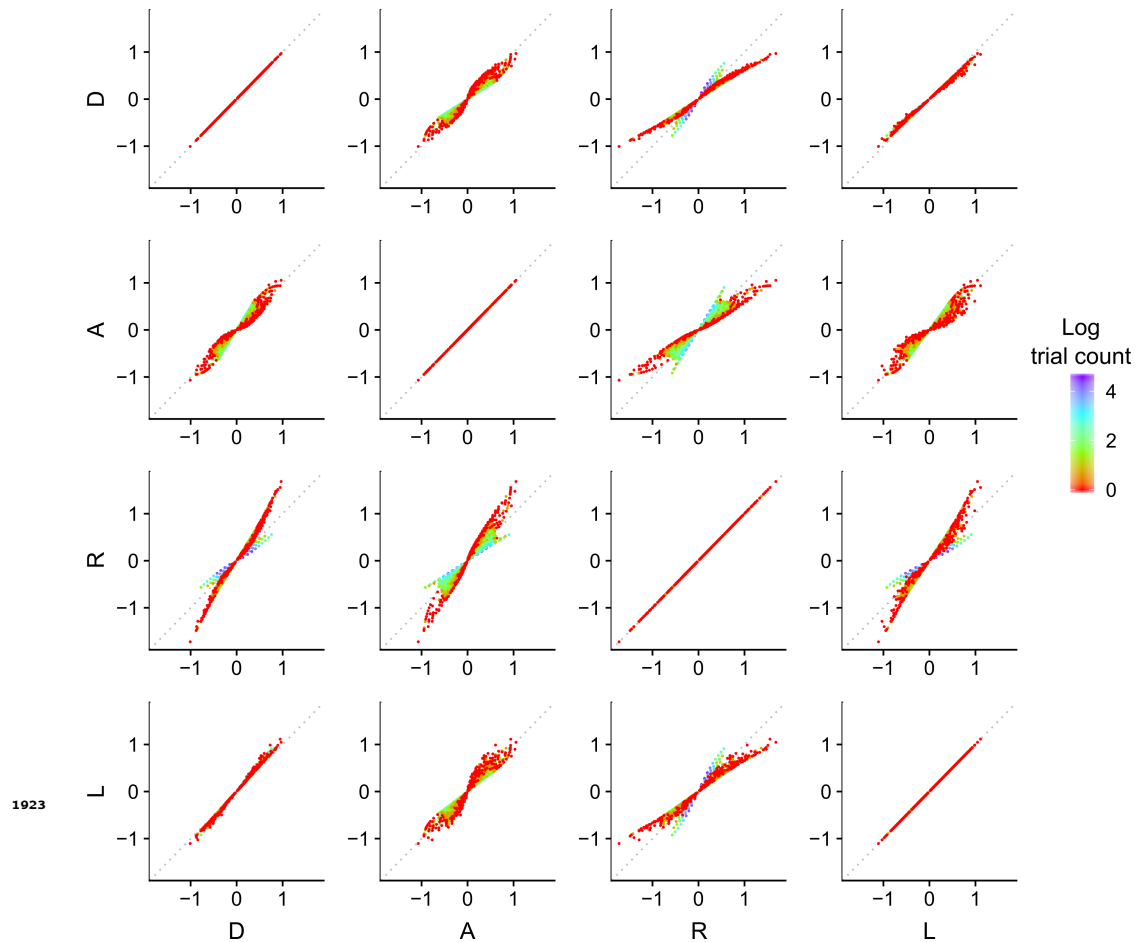


Figure 3-Figure supplement 2. Discriminability differences between the different encoding rules. This figure illustrates the discriminability differences between the different encoding rules considered in this study. Each dot represents the discriminability value for a pair of numerosity values v_1 and v_2 presented on a given trial to the participants in Experiment 1. For the sampling models, the discriminability rule is defined as

$$\frac{\theta(v_1) - \theta(v_2)}{\sqrt{\theta(v_1)(1 - \theta(v_1)) + \theta(v_2)(1 - \theta(v_2))}},$$

where θ corresponds to the respective Accuracy maximizing (A), Reward maximizing (R) or Decision by Sampling (D) encoding rules. For the logarithmic model (L) the discriminability rule is defined as

$$\log(v_1) - \log(v_2).$$

The color of each dot represents the log of the number of occurrences for the pairs of input values v_1 and v_2 . Note that the encoding values of the presented numerosities are different depending on the encoding rule, which makes it possible to identify the participants' encoding strategy. Also note that for our imposed prior distribution, the DbS encoding rule is similar to the logarithmic rule, which explains the smaller difference in the quantitative predictions between these two models. Nevertheless, DbS was always the model that provided the best quantitative and qualitative predictions irrespective of incentivized goals.

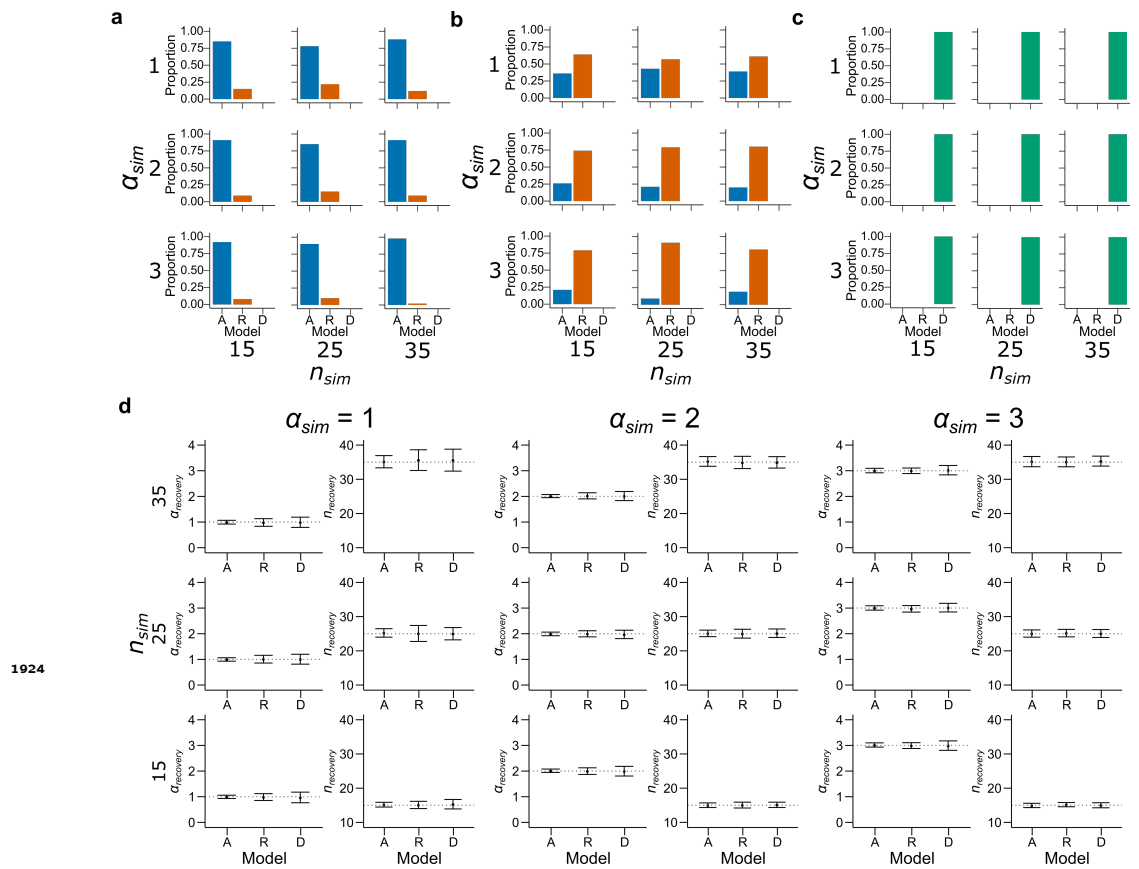


Figure 3-Figure supplement 3. Model recovery with both α and n as free parameters. Synthetic data preserving the trial set statistics and number of trials per participant used in Experiment 1 was generated 100 times for each encoding rule with various values of α and n . A model for each encoding rule was fit to the data using maximum likelihood estimators with α and n as free parameters. The histograms represent the proportion best fitting models for **a**) Accuracy, **b**) Reward and **c**) DbS models. Results are shown for different values of α (top: $\alpha = 1$, middle: $\alpha = 2$ and bottom: $\alpha = 3$) and n (left: $n = 15$, middle: $n = 25$ and right: $n = 35$). While DbS is always well recovered, the Accuracy and Reward models tend to be confounded with each other. **d**) This same synthetic data was fit with its generating model with α and n as free parameters using maximum likelihood estimators. Results are shown for different values of α (first and second columns: $\alpha = 1$, third and fourth columns: $\alpha = 2$ and fifth and sixth columns: $\alpha = 3$) and n (top: $n = 35$, middle: $n = 25$ and right: $n = 15$). Error bars represent one standard deviation of the recovered parameter across simulations. The parameters are well recovered by the respective generating model.

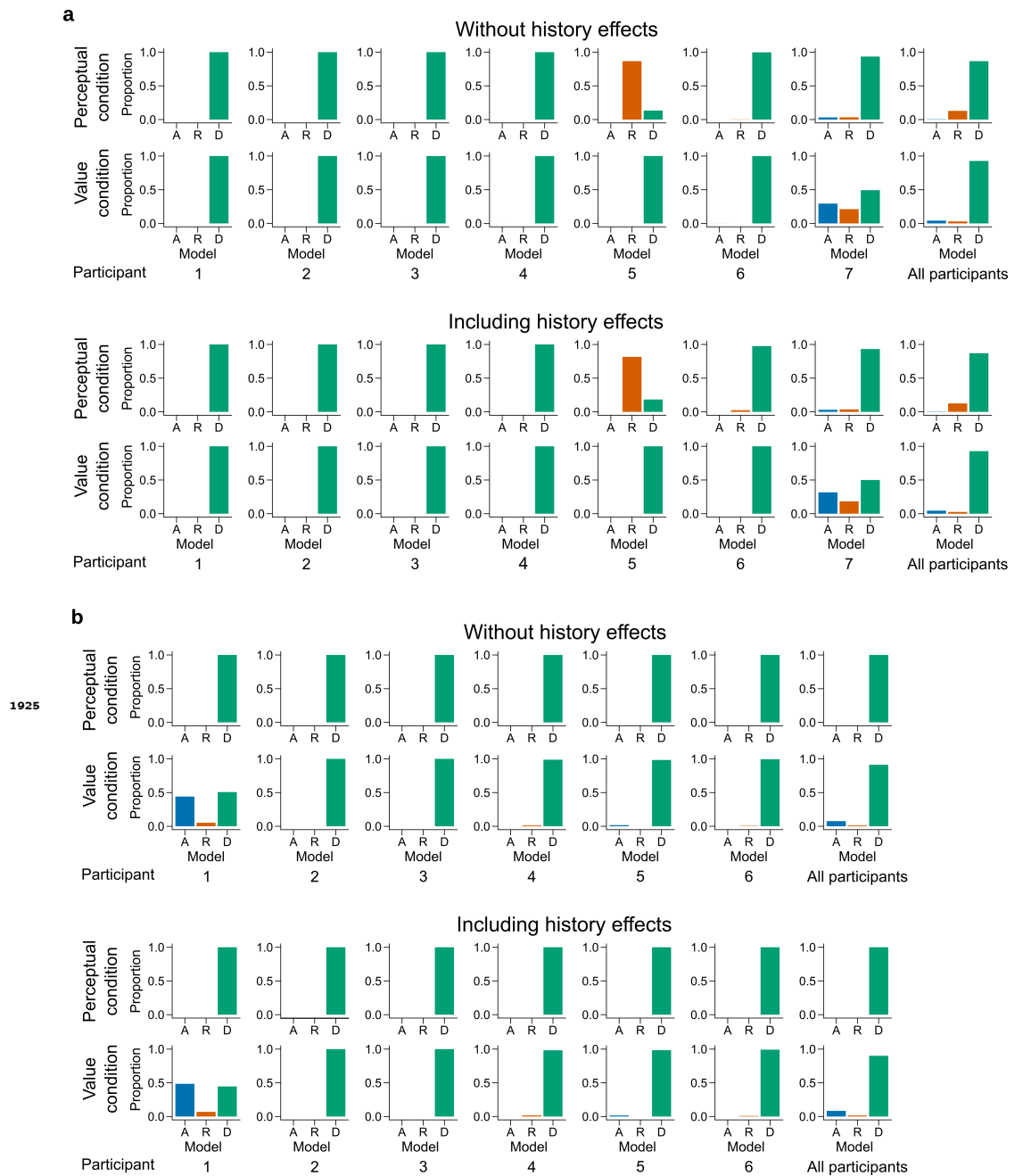


Figure 4-Figure supplement 1. Latent mixture model fits for each participant. Individual level fit of the latent mixture model excluding (*top*) or including (*bottom*) choice history effects for **a**) Experiment 1 and **b**) Experiment 2. The panels on the far right shows the average fit for all the participants of the given experiment. DbS is strongly favored for nearly all participants and clearly favored across participants, irrespective of the experimental condition. Including choice and correctness information of previous trials has minimal influence in the results of these analyses, which rules out the influence of these effects on the decision rule used by the participants.

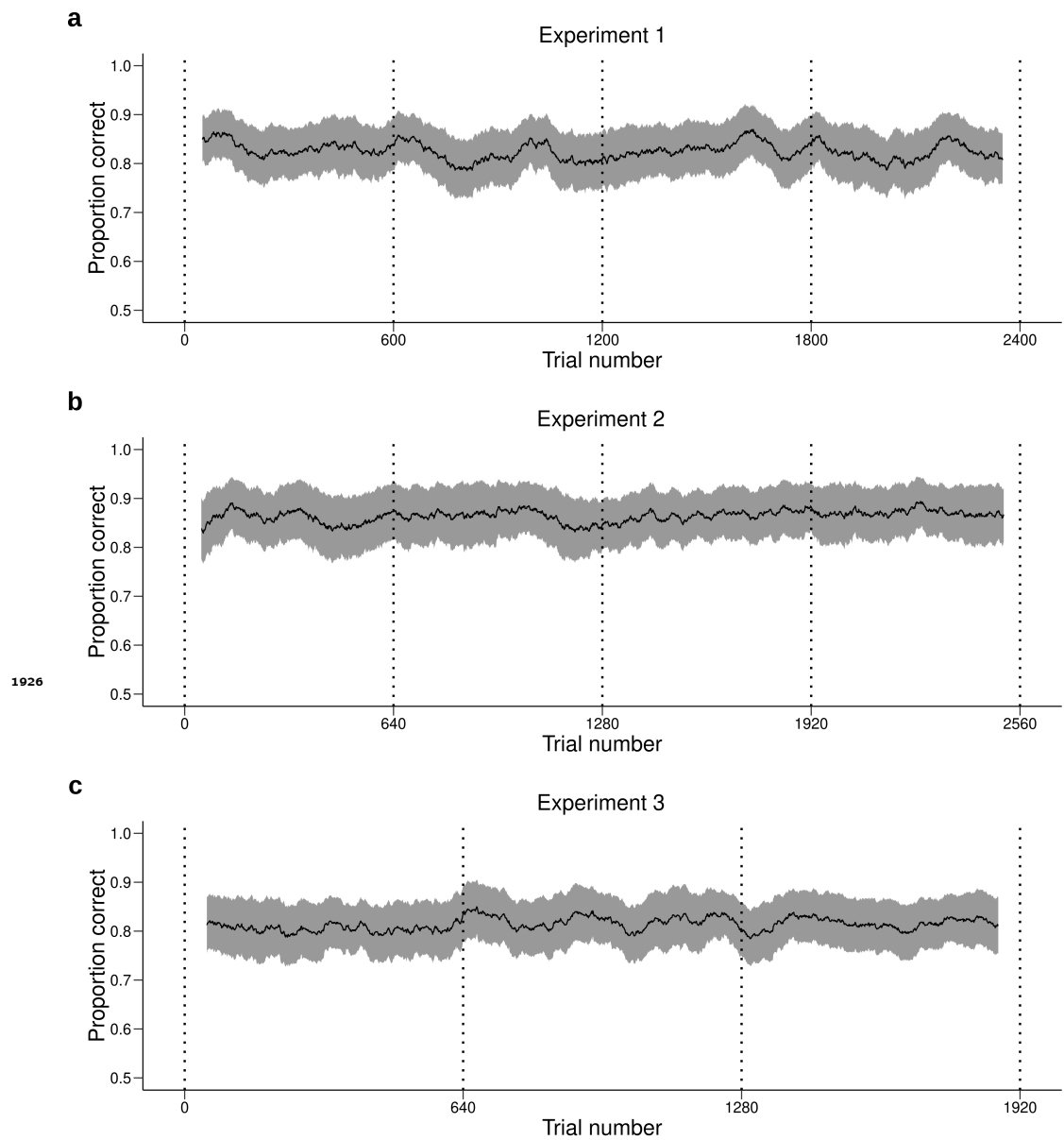


Figure 4-Figure supplement 2. Performance across time. Behavioral performance (mean \pm s.e.m. across participants) averaged over a moving window of 100 trials for **a)** Experiment 1, **b)** Experiment 2 and **c)** Experiment 3. Each daily session took place between two dotted vertical lines. The performance of the participants is stable during and between daily sessions. Therefore, the quantitative and qualitative results presented in the main text are not likely to be influenced by changes in performance over time.

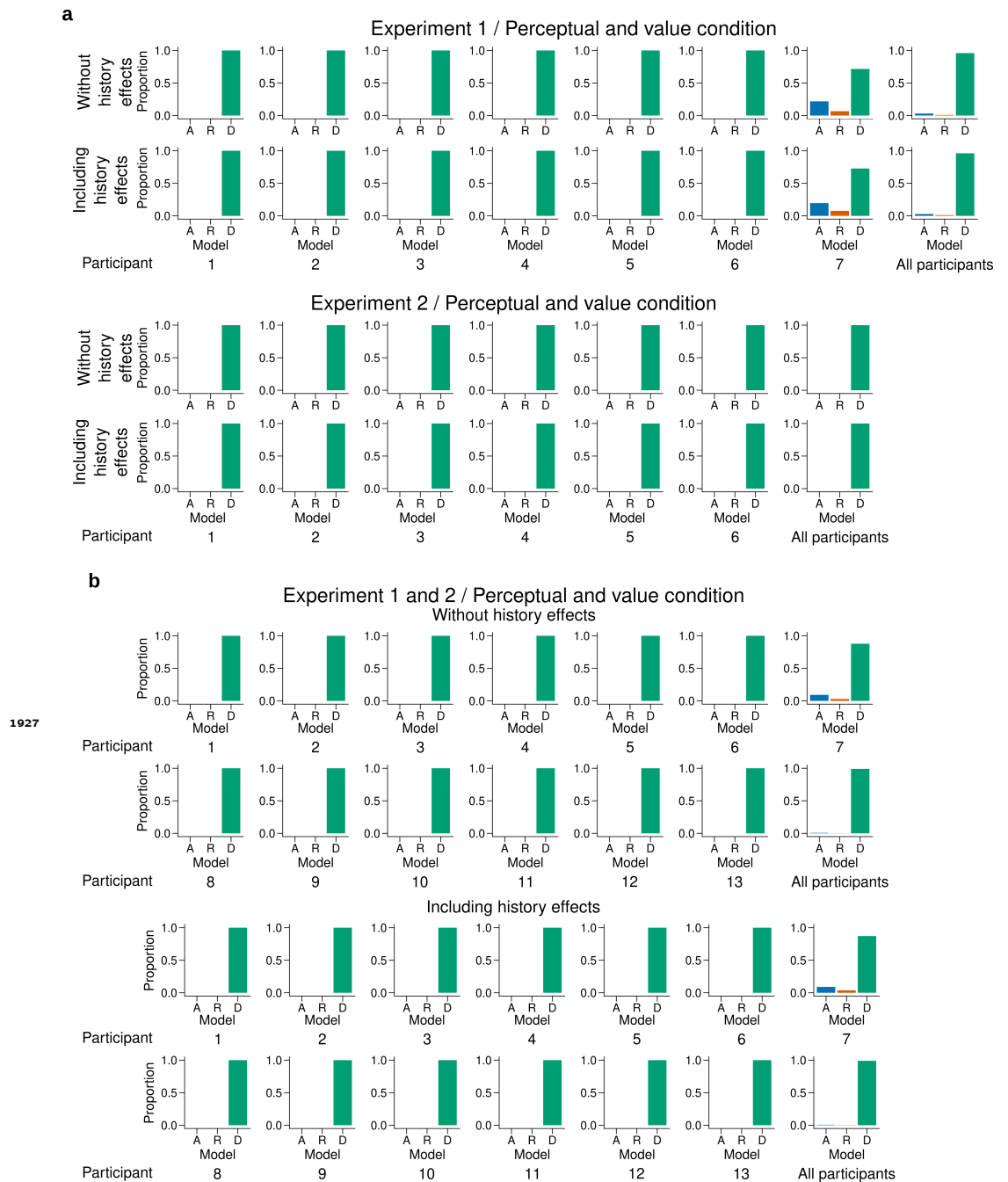
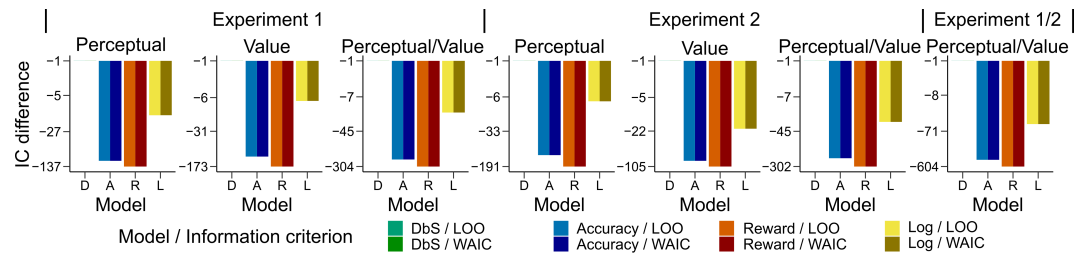
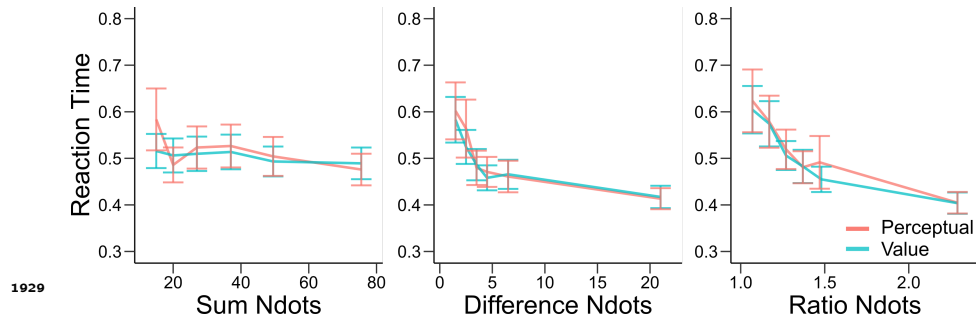


Figure 4-Figure supplement 3. Individual level fit of the latent mixture model combining data across experiments and experimental conditions. Individual level fit of the latent mixture model combining data across both experimental conditions for Experiment 1 (*top*) and Experiment 2 (*bottom*). **b)** Individual level fit of the latent mixture model combining data across both experimental conditions and both experiments. Each panel shows the results excluding (*top*) or including (*bottom*) choice history effects. The panels labeled "All participants" show the average fit for all the participants of the given experiment. DbS is strongly favored irrespective of incentivized goals. Including the previous trial effects has minimal influence on these results.



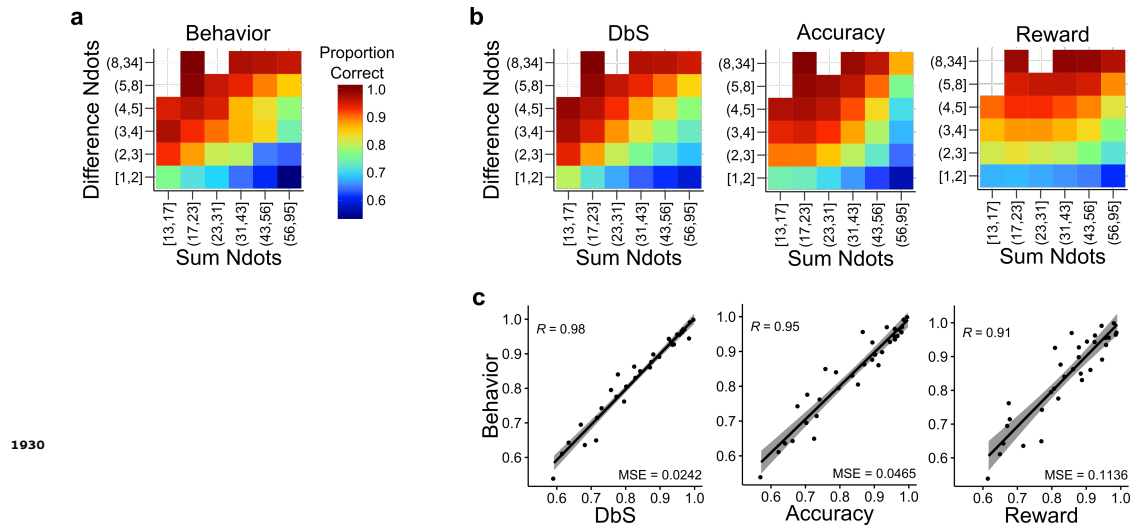
1928

Figure 4-Figure supplement 4. Model comparison based on leave-one-out cross-validation metrics. Quantitative comparison of the models including choice and correctness effects of previous trials based on leave-one-out cross-validation metrics. **a)** Difference in LOO and WAIC between the best model (DbS (D) in all cases) and the competing models: Accuracy (A), Reward (R) and Logarithmic (L) models. Each panel shows the data grouped for each and across experiments and experimental conditions (see titles on top of each panel). Including the previous choice and correctness effects has only little influence on the results (compare with Figure 4b in main text). The DbS model provides the best fit to the behavioral data.



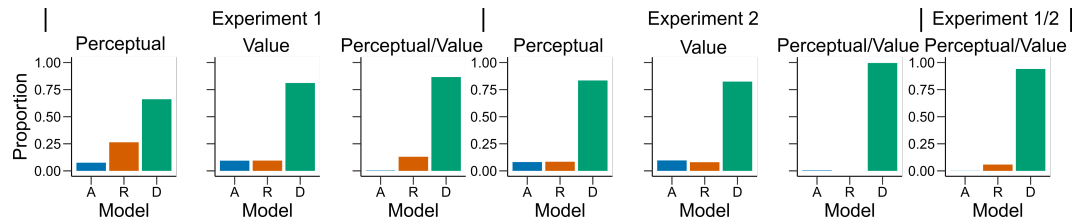
1929

Figure 4-Figure supplement 5. Reaction times are similar in the perceptual and value conditions. Mean reaction times of participants in experiments 1 and 2 in the perceptual (red) and value (blue) condition. Error bars represent s.e.m. across participants. Reaction times are presented as a function of the sum of the number of dots in both clouds (left), the absolute difference between the number of dots in both clouds (middle) and the ratio of the number of dots in the most numerous cloud over the less numerous cloud (right). Non-parametric ANOVA tests revealed no significant differences in any of these behavioral assessments (all tests $P > 0.4$).



1930

Figure 4-Figure supplement 6. Behavior and model predictions as a function of sum and difference in dots. **a)** Average behavior in both conditions of experiments 1 and 2 as a function of the sum of the number of dots in both clouds (Sum Ndots) and the absolute difference between the number of dots in both clouds (Difference Ndots). The data is binned as in Figure 4 but now expanded in two dimensions. **b)** Predictions of each encoding rule model fit with only n as a free parameter shown with the same scale as in a. **c)** Linear regression between the behavior for each combination of Sum Ndots and Difference Ndots bins and the predictions of each model for the same bins. DbS captures best the changes in behavior across bins of sum and absolute difference of the number of dots in both clouds. This analysis should not be considered as a quantitative proof, but as a qualitative inspection of the results presented in Figure 4.



1931

Figure 4-Figure supplement 7. Model fit for the first experimental condition of each participant. Similar as in Figure 4a., bars represent proportion of times an encoding rule (Accuracy (A, blue), Reward (R, red), DbS (D, green)) was selected by the Bayesian latent-mixture model based on the posterior estimates across participants. Each panel shows the data grouped for each and across experimental conditions and experiments (see titles on top of each panel). The latent-mixture model was only fit to the first condition that was carried out by each participant. As the participants did not know of the second condition before carrying it out, they could not adopt compromise strategies between the two objectives. Therefore, the fact that DbS is favored in the results is not an artifact of carrying out two different conditions in the same participants.

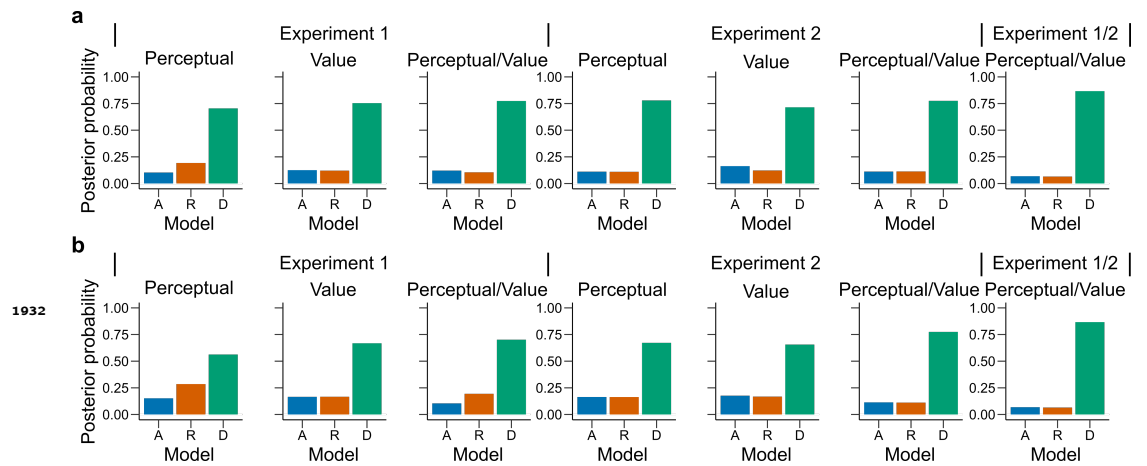


Figure 4-Figure supplement 8. Latent vector π posterior estimates. Bars represent the posterior distribution of the latent vector π , with each bar representing an encoding rule (Accuracy (A, blue), Reward (R, red), DbS (D, green)). Results are presented for all **(a)** sessions and **(b)** only the first condition carried out by each participant. DbS is consistently the most likely encoding rule.

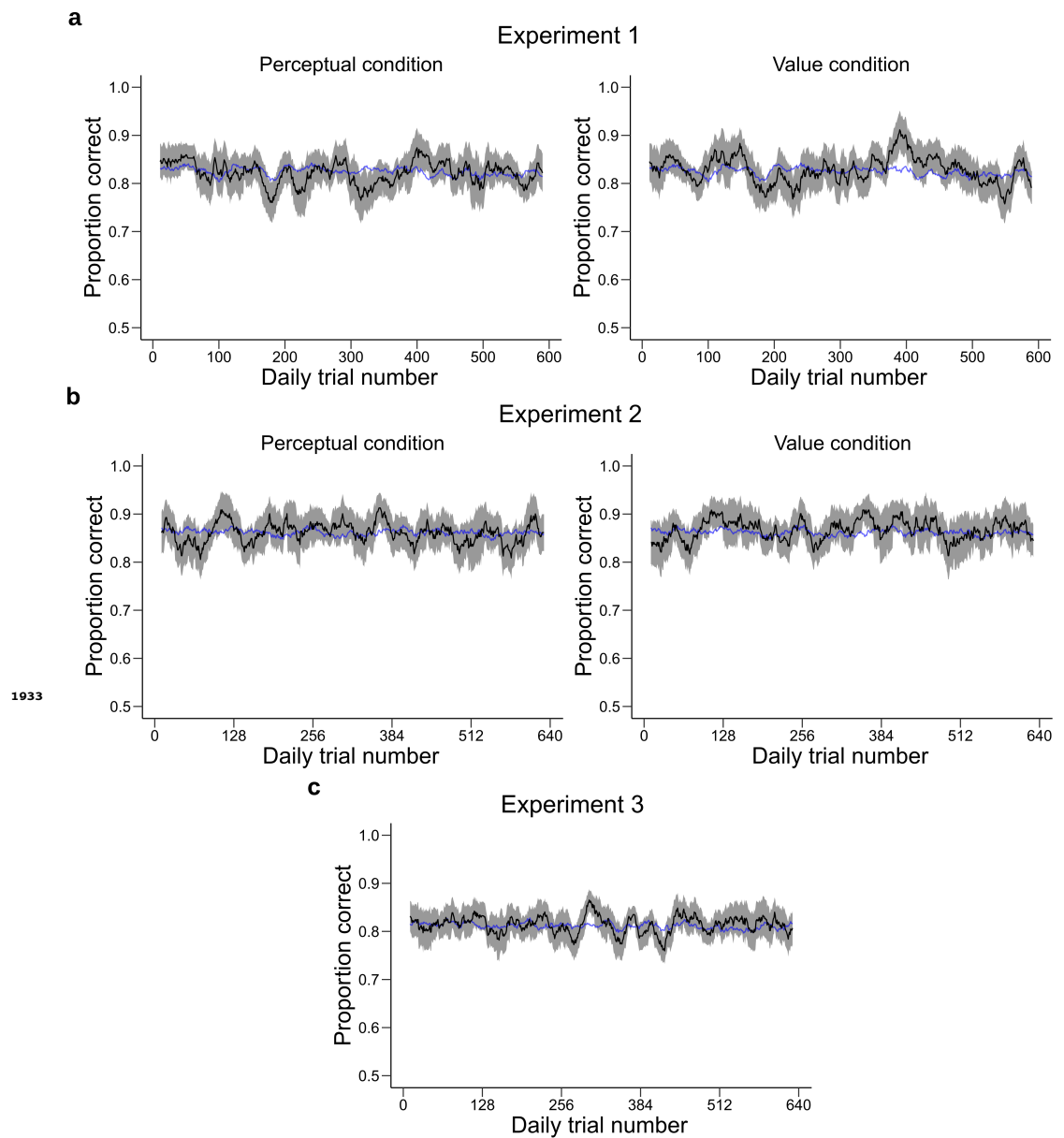
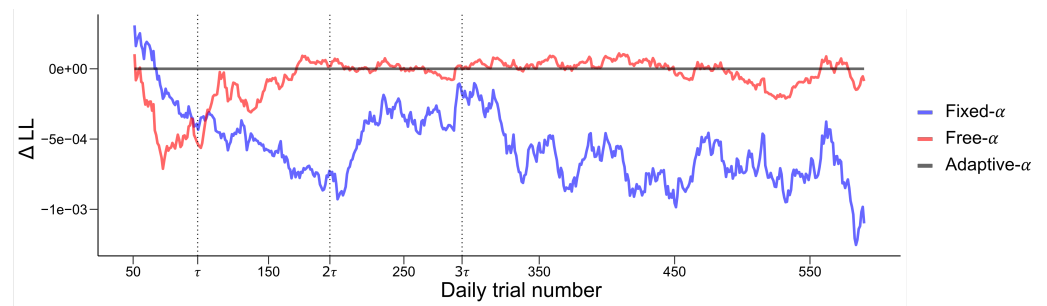


Figure 5-Figure supplement 1. Performance across trial experience. These plots represent the performance of the participants as a function of the number of trials they have experienced during the session. The performance of the participants (black, shaded area represents \pm s.e.m. across participants) was averaged over a moving window of 21 trials and is shown for experiment 1 (**a**) experiment 2 (**b**) and experiment 3 (**c**). The blue line represents the performance predicted by the α -adaptation model using the same moving window average. The model provides a good fit to average performance.



1934 **Figure 5-Figure supplement 2.** Quantitative and dynamical analysis of adaptation over time. To further investigate the adaptation of the prior, we fit three model of varying complexity to the data of experiments 1, 2 and 3. The Fixed- α model (blue) is defined with a fixed $\alpha = 2$. The Free- α model (red) allows the α parameter to vary across participants but is kept constant across time. The Adaptive- α corresponds to the model presented in Figure 5 where the prior adapts as the participants gains experience with the experimental distribution of dots. To allow a fair comparison with the Free- α model, the δ parameter, corresponding to the asymptotic value of the prior, was free to vary across participants. The log-likelihood of each model on each trial were averaged over a moving window of 100 trials and the log-likelihood of the Adaptive- α model was subtracted for comparison. Vertical dashed lines represent 1, 2 and 3 times τ , where τ controls the rate of adaptation in the Adaptive- α model. The Adaptive- α model provides a better fit for the first trials (until around 2τ), these trials correspond to the adaptation period where the α parameter is changing in the Adaptive- α model (see **Figure 5**). After this point the Adaptive- α and Free- α models provide a similar fit. This is to be expected as the function controlling the decay of α reaches its asymptotic value, leaving the two model virtually identical. The Fixed- α provides overall a worse fit, except for the early trials.