

## RESEARCH ARTICLE

## Efficient numerosity estimation under limited time

Joseph A. Heng <sup>1,2</sup>, Michael Woodford<sup>3</sup>, Rafael Polania <sup>1,2\*</sup>

**1** Decision Neuroscience Lab, Department of Health Sciences and Technology, ETH Zurich, Zurich, Switzerland, **2** Neuroscience Center Zurich, Zurich, Switzerland, **3** Department of Economics, Columbia University, New York, New York, United States of America

\* [rafael.polania@hest.ethz.ch](mailto:rafael.polania@hest.ethz.ch)

 OPEN ACCESS

**Citation:** Heng JA, Woodford M, Polania R (2025) Efficient numerosity estimation under limited time. *PLoS Comput Biol* 21(3): e1012790. <https://doi.org/10.1371/journal.pcbi.1012790>

**Editor:** Tobias U. Hauser, University of Tübingen: Eberhard Karls Universität Tübingen, GERMANY

**Received:** October 13, 2023

**Accepted:** January 13, 2025

**Published:** March 7, 2025

**Copyright:** © 2025 Heng et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data availability statement:** This work analyzes the data from a previously published dataset available on the Open Science Foundation at <https://osf.io/svcy5/>. The code produced in the work is available on the ETH Zurich Research Collection at <https://doi.org/10.5905/ethz-1007-853>.

**Funding:** This work was supported by a European Research Council (ERC) starting grant (ENTRAINER) to R.P. This project has

## Abstract

The ability to rapidly estimate non-symbolic numerical quantities is a well-conserved sense across species with clear evolutionary advantages. However, despite its importance, this sense is surprisingly imprecise and biased, and a formal explanation for this seemingly irrational behavior remains unclear. We develop a unified normative theory of numerosity estimation that parsimoniously incorporates in a single framework information processing constraints alongside (i) Brownian diffusion noise to capture the effects of time exposure of sensory information, (ii) logarithmic encoding of numerosity representations, and (iii) optimal inference via Bayesian decoding. We show that for a given allowable biological capacity constraint our model naturally endogenizes time perception during noisy efficient encoding to predict the complete posterior distribution of numerosity estimates. This model accurately predicts many features of human numerosity estimation as a function of temporal exposure, indicating that humans can rapidly and efficiently sample numerosity information over time. Additionally, we demonstrate how our model fundamentally differs from a thermodynamically-inspired formalization of bounded rationality, where information processing is modeled as acting to shift away from default states. The mechanism we propose is the likely origin of a variety of numerical cognition patterns observed in humans and other animals.

## Author summary

Humans can estimate the number of elements in a set without counting. We share this ability with other species, suggesting that it is evolutionarily relevant. However, this sense is variable and biased. What is the origin of these imprecisions? We take the view that they are the result of an optimal use of limited neural resources and limited processing time. Because of these limitations, stimuli are encoded with noise. The observer then optimally decodes these noisy representations, taking into account its knowledge of the distribution of stimuli. We build on this perspective by incorporating stimulus presentation time directly into the encoding process. This model can parsimoniously predict key

received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No. 758604) to R.P., and support from the U.S. National Science Foundation, under grant SES-DRMS-1949418 to M.W. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing interests:** The authors have declared that no competing interests exist.

characteristics of our perception and outperforms quantitatively and qualitatively a popular modeling approach that considers resource limitations at the stage of the response rather than the encoding.

## Introduction

The ability to rapidly represent and estimate non-symbolic numerical quantities is a fundamental cognitive function for behavior in humans and other animals, which may have emerged during evolution to support fitness maximization [1]. Since the properties of numerosity estimation started to be studied nearly a century ago, it has been commonly observed that the representation and estimation of numerical quantities are imprecise and biased [2]. Despite the importance of numerosity estimation for various cognitive processes and ultimately survival, the questions remain: what are the origins of the observed variability and biases in numerosity estimations? Are these deviations efficient and predictable when organisms are urged to rapidly estimate numerical quantities?

Extensive empirical research in the representation and estimation of non-symbolic numerical quantities has consistently reported and studied various features that characteristically emerge during numerosity estimation, including: (i) subitizing small numbers [3]; (ii) overestimation of small numbers (outside the subitization range) and underestimation of large numbers [4], with especially biased estimates in the case of larger numbers [5]; (iii) a coefficient of variation that is approximately constant across all numerosities, a property termed scalar variability [6]; and (iv) estimation acuity modulated by duration of stimulus presentation and sensory reliability [7]. But do all the above-mentioned behavioral patterns have a common origin?

The last decades have been marked by the development of models of behavior in which perception has been proposed to be instantiated as a Bayesian inference process. This suggests that our nervous system jointly considers the environmental (or contextual) distribution of sensory stimuli and the unreliability of the signals perceived by the observer. This approach has been instrumental in explaining in a parsimonious manner a variety of behavioral biases including underestimation, overestimation, and the degree of variability of estimated magnitudes and quantities [8]. However, this approach does not explicitly consider the different sets of constraints that biological systems face when interacting with the environment. This is a fundamental aspect to consider in any formulation that attempts to explain the behavior of biological systems given the fact that organisms do not have unlimited biological resources or unlimited time to process sensory information from the environment, and moreover, neural computations are metabolically expensive [9]. Thus, it has been suggested that the observed variability and biases in our estimations of our sensory world emerge from efficient processes based on fundamental principles of encoding information from environments with statistical regularities [10–14].

Here we argue that many of the above-mentioned behavioral features emerging during numerosity estimation have a common origin: given biological constraints on information acquisition, numerosity estimation emerges from a system that efficiently considers, first, prior knowledge of the environment, second, information of the current numerosity being evaluated, and third, the amount of time available to process such information.

We develop a unified normative model of numerosity estimation that parsimoniously incorporates information constraints together with long modeling traditions of human and animal psychophysical performance in psychology and neuroscience: (i) Brownian diffusion noise to capture the effects of time exposure of sensory information [15], (ii) logarithmic

encoding of numerosity representations [16], and (iii) optimal Bayesian decoding. As a result, we show that for a given allowable biological capacity constraint, our model naturally incorporates time perception during noisy efficient encoding to predict the corresponding posterior distribution of numerosity estimates via optimal Bayesian decoding. Here we refer to our approach as the “sequential-encoding/Bayesian-decoding” model, henceforth SEB.

We also consider a second well-known approach for studying bounded rationality inspired by principles of thermodynamics and statistical physics. This family of models assumes that given a default state (e.g., a default distribution over possible responses) and a sensory stimulus, the observer acts in a way such that they attempt to shift from the default state to a new state that matches as closely as possible the value of the sensory stimulus. Bounded rationality comes into play in the case of acting when only a given amount of change in information (energy invested) between the default and new state can be afforded. This class of models has been used in a wide range of applications [17–20], including recently to study how perceptual estimation under limited time relates to cognitive capacity and action responses [21]. Here we refer to this class of models as the “thermodynamically inspired model”, henceforth TIM.

A key contribution of our work is the formal demonstration that the two approaches that we consider here (SEB and TIM) are in fact classes of models with completely different views on bounded rationality. To avoid confusing them, we clarify their differences here. On the one hand, variability in the estimation responses in SEB is attributed to *sensing* costs, which generate noisy sensory encoding. On the other hand, in instantiations of TIM applied to sensory estimation, variability is generated by *acting* costs during response selection. Crucially, here we demonstrate that these two approaches applied to numerosity estimation lead to apparently similar but distinguishable quantitative and qualitative predictions that are identifiable and falsifiable. Our empirical tests applied to a large numerosity estimation data set provide a clear indication that humans follow a SEB rather than a TIM approach, meaning they can rapidly and efficiently sample numerosity information over time via an efficient noisy encoding and Bayesian decoding process.

## Results

The presentation of our results is divided into three parts: First, we present our sequential-encoding/Bayesian-decoding model (SEB) which parsimoniously endogenizes perceptual exposure times in its likelihood function alongside parameters of the prior distribution for a given biological capacity bound. Second, we introduce the thermodynamically inspired model (TIM) applied to sensory estimation, and compare it with the SEB model. Third, we apply rigorous quantitative and qualitative model evaluations based on a large publicly available human numerosity estimation dataset ( $n = 400$  participants across four different experiments).

### A Bayesian model of numerosity estimation

Extensive behavioral and physiological work studying the representation of both non-symbolic and symbolic numerical quantities strongly suggests that internal representations  $r$  can be assumed to be encoded by a quantity that is proportional to the logarithm of the number  $n$  plus stimulus-independent random error that is assumed to be normally distributed and unbiased [5,16,22]

$$r \sim N(\log n, v^2). \quad (1)$$

However, a key contribution of our work is to formally study how these perceptual errors may depend on stimulus duration  $t$  of the form

$$r \sim N(m(n), v^2(t)), \quad (2)$$

which represents the likelihood in our Bayesian framework. Here,  $m(n)$  is an affine function of  $\log(n)$ . Below, we will develop the theory by finding the parameters of the encoding process that minimize the MSE between the inputs and the estimation. That is, we aim to optimize the encoding function based on the variability of the representation which is a function of the sensory representation as a function of stimulus duration  $t$ .

We assume the prior distribution to be a log-normal distribution from which the true numerosity  $n$  is drawn to be

$$\log n \sim N(\mu, \sigma^2). \quad (3)$$

Note that  $\mu$  and  $\sigma$  are the expected value and standard deviation of the random variable's natural logarithm, and not the expectation and standard deviation of  $n$  itself.

While the distribution of various quantities in linguistics, economics, and ecology appears to be well-described by log-normal distributions [23], others have argued that power-law distributions approximately describe the empirical frequency of numbers in natural environments [24,25]. We note, however, that the two-parameter family of possible log-normal prior distributions includes as a limiting case the power-law distributions (see S1 Note for proof). In brief, we consider a normalized prior of the form

$$p(n) \propto \exp(-\alpha(\log n) - \gamma(\log n)^2), \quad (4)$$

for some parameters  $\alpha, \gamma$  with  $\gamma \geq 0$ . If  $\gamma > 0$ , this corresponds to a log-normal prior, with  $\mu = (1 - \alpha)/(2\gamma)$ ,  $\sigma^2 = 1/(2\gamma)$ . If instead  $\gamma = 0$  but  $\alpha > 0$ , this corresponds to a power-law prior

$$p(n) \propto n^{-\alpha}. \quad (5)$$

Thus, our model allows for the possibility that encoding and decoding are adapted to different priors that are learned for different contexts, rather than a single process being used in all contexts. However, in the following theoretical developments, we consider a log-normal prior for simplicity.

Here we assume that the objective of the decision-maker is to obtain numerosity estimates  $\hat{n}$  that minimize the MSE when stimuli are drawn from the prior distribution. It can be shown that this implies that conditional on  $n$ , the estimate  $\hat{n}$  will be log-normally distributed (S1 Note)

$$\log \hat{n} | n \sim N(\hat{\mu}(n, t), \hat{\sigma}^2(t)), \quad (6)$$

where  $\hat{\mu}(n, t)$  is an affine function of  $\log n$ , and  $\hat{\sigma}^2(t)$  is independent of  $n$ . However, both  $\hat{\mu}$  and  $\hat{\sigma}^2$  may depend on temporal numerosity processing  $t$ , as we formally elaborate below.

**Exposure time and precision of internal representations.** A key hallmark in the development of our theoretical framework is that we now assume that the sensory evidence of the input stimulus is given by a Brownian motion with a drift that depends on the stimulus (formally defined below). Thus, by modeling sensory percepts in this way, we follow a

long modeling tradition of process models of perception and action that includes the popular drift-diffusion model (DDM) [15]. Models of this kind have been used since the late 60s to account quantitatively for the way in which the accuracy of perceptual judgments is affected by manipulations of viewing time [26].

Formally, we now suppose that the internal representation  $r$  consists of the sample path of a Brownian motion  $z_s$  over a time interval  $0 \leq s \leq \tau$ , starting from an initial value  $z_0 = 0$ . The drift-diffusion parameter  $m$  of the Brownian motion is assumed to depend on  $n$ , while its instantaneous variance  $\omega^2$  is independent of  $n$ ; the length of time  $\tau$  for which the Brownian motion evolves is also independent of  $n$ , but depends on the viewing time  $t$ . In other words, we assume that the agent makes observations of momentary evidence  $\delta_{z_i} \sim N(m(n)\delta_t, \omega^2\delta_t)$  in small steps  $i$  for an infinitesimal duration  $\delta_t$ , where the accumulated evidence at step  $k$  is given by  $z_s(k\delta_t) = \sum_{i=1}^k \delta_{z_i}$ .

Please note that the instantaneous variance  $\omega^2$  of the diffusion process should not be confused with the resulting encoding noise  $v^2$  in the logarithmic space. In fact, our goal will be to formally find how  $v$  depends on elements that determine the diffusion process  $m$  and  $\omega$  for given stimulus  $n$  and viewing time  $t$ .

More specifically, under the assumption that the particle position under Brownian motion is normally distributed with its parameters evolving as a function of  $\tau$ , one can show that  $r$  is a draw from the distribution (S2 Note for details)

$$r \sim N(m(n), \omega^2/\tau). \quad (7)$$

Our goal is now to find a solution of how such a dynamic perceptual system should operate under limited resources. Crucially, we suppose the average value of  $m^2$  is subject to a power constraint, that is to be within some finite bound

$$E[m^2] \leq \Omega^2 < \infty. \quad (8)$$

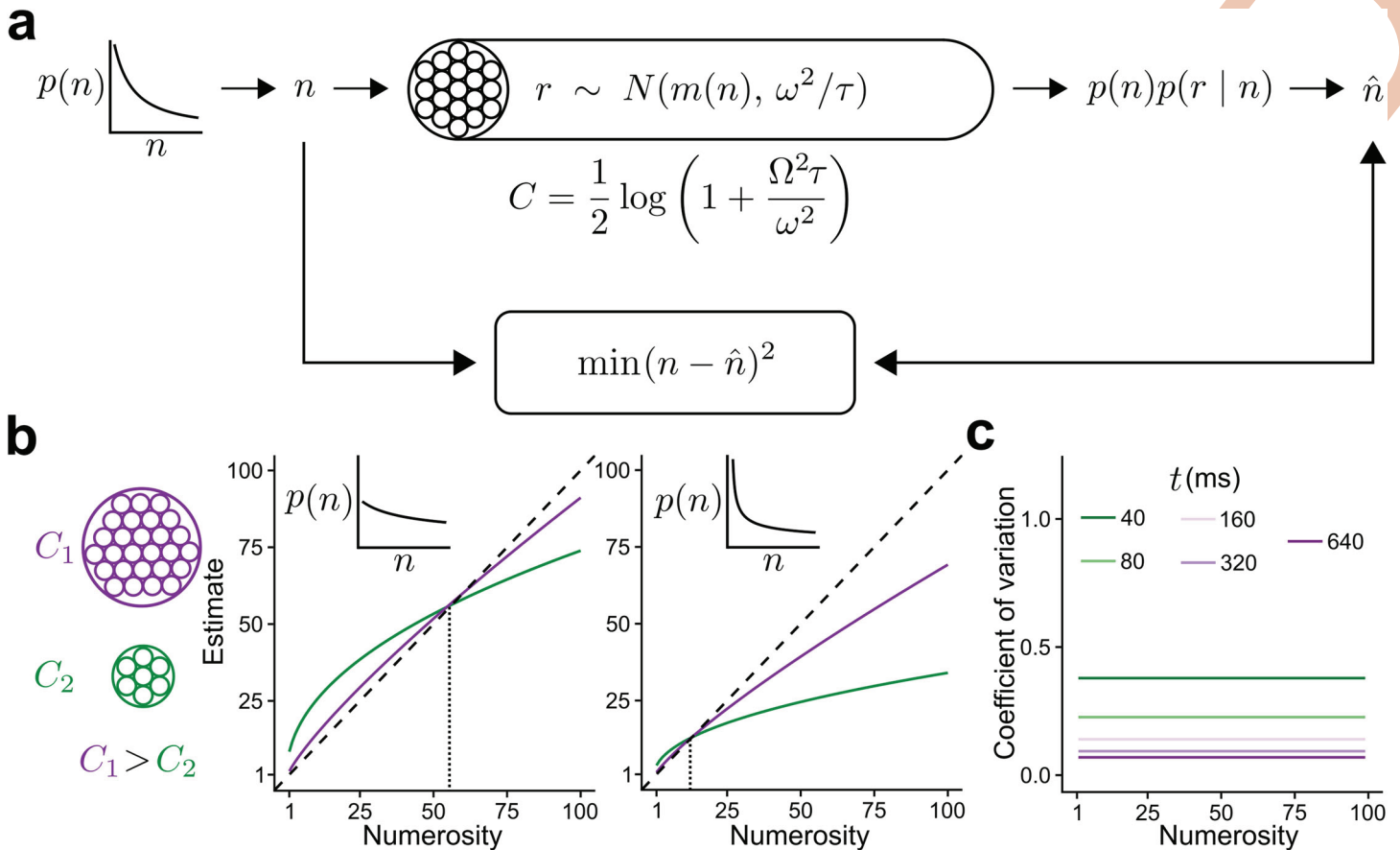
This bound on the amount of variation in the drift limits the precision with which different stimuli can be perceived for any given  $\tau$ . The value of  $\tau$  is assumed to grow linearly with the viewing time, up to some time bound  $t_{max}$ ,

$$\tau = \min(t, t_{max}), \quad (9)$$

representing a constraint on the amount of time that the decision maker is willing to invest in the accumulation of evidence. The latter bound constrains the degree to which precision can be increased by further increases in viewing time.

The definition of this optimization problem with constraints effectively states that  $r$  can be seen as the output of a *Gaussian channel* with input  $m$  [27] that depends on the input stimulus  $n$ ; hence the problem of optimally choosing the function  $m(n)$  is equivalent to an optimal encoding problem for a Gaussian channel (S2 Note). The capacity  $C$  of such a channel is a quantitative upper bound on the amount of information that can be transmitted regardless of the encoding rule, which is equal to (S3 Note)

$$C = \frac{1}{2} \log \left( 1 + \frac{\Omega^2 \tau}{\omega^2} \right). \quad (10)$$



**Fig 1. Overview of the SEB model.** **a)** Schematic description of the SEB model. A numerosity  $n$  is drawn from a prior distribution  $p(n)$ . The observer has a limited capacity  $C$  to represent the numerosity. The internal representation  $r$  is a random draw from a Gaussian distribution, the mean of which depends on  $n$  but the variance does not. The observer then infers the estimate  $\hat{n}$  based on the representation  $r$  and the prior distribution of  $n$  as to minimize the MSE between the estimate and the numerosity. **b)** Illustration of the predictions for an observer with a high (purple) or low (green) channel capacity where  $n$  is drawn from a distribution with a high (left) or low (right) variance. All curves exhibit overestimation for lower numerosities and underestimation for higher numerosities. However, these biases are reduced in the case of high capacity. The crossover point between under- and overestimation increases with the variance of the numerosity distribution. **c)** Illustration of the coefficient of variation (i.e.,  $SD[\hat{n}]/E[\hat{n}]$ ) for different capacities. The coefficient of variation is independent of the numerosity and decreases with capacity, which is dependent on the viewing time  $t$ .

<https://doi.org/10.1371/journal.pcbi.1012790.g001>

Note that in our model the channel capacity  $C$  grows as a logarithmic function of  $\tau$  because the correlation of successive increments in the encoding by a Brownian motion prevents the information content from growing linearly in proportion to such increments.

Here we assume that the goal is to design a capacity-limited system that minimizes the mean squared error (MSE) of the estimate  $\hat{n}$  when  $n$  is drawn from a log-normal prior distribution (i.e., the same objective function stated in the previous section, see Fig 1). It is possible to show that in our optimization problem, which assumes a channel with “power transmission” constraint  $\Omega^2$ , the optimal drift function is

$$m = \xi + \psi \log n \tag{11}$$

with

$$\xi = -\psi\mu, \quad \psi = \frac{\Omega}{\sigma}, \tag{12}$$

and the encoding noise  $\nu$  in Eq 2 is given by (see S2 Note for proof)

$$\nu(t) = \frac{1}{\sqrt{\tau}} \frac{\omega}{\Omega} \cdot \sigma. \quad (13)$$

That is, encoding precision grows with viewing time  $t$  (until reaching the bound  $t_{max}$  if the stimulus is presented long enough). Recall that  $\sigma$  is the variance of the log-normal prior, and therefore the solution reveals that the likelihood is independent of parameter  $\mu$  of the log-normal prior distribution, but depends on the second moment of this prior distribution and viewing time  $t$ . Defining  $R \equiv \Omega/\omega$ , the noise of numerosity encoding is given by  $\nu(t) = 1/G$ , where  $G = \min(R\sqrt{t}/\sigma, B)$  and  $B$  a maximum biologically allowed bound on sensory precision related to  $t_{max}$  ( $B = R\sqrt{t_{max}}/\sigma$ ). Note that the two information bounds affect the model differently. The higher is  $\Omega$ , the more information gathered per time unit, whereas  $B$  captures the maximum amount of information that can be gathered. Notice that this solution implies a multiplicative trade-off between  $\Omega$ ,  $\omega$ , and  $t$  (similar to the standard DDM). However, this relation may look different under other assumptions (e.g., non-uniform noise in the encoding space).

These results lead to the following predictions from our model: (i)  $E[\hat{n} | n]$  is a concave function of  $n$  with overestimation for small numbers (when these are not so small that the discreteness of available responses leads to nearly-deterministic responses), but underestimation for large numbers (Fig 1b and S2 Note). (ii) The crossover point from overestimation to underestimation changes as a function of the numerosity range and prior variance (see S4 Note). In addition, the concavity of  $E[\hat{n} | n]$  depends on the amount of resources available to perform the numerosity estimation task. This prediction was clearly confirmed in a previous empirical work [4]. (iii) Because of the discreteness in the set of responses, there is predicted to be little variability in responses in the case of small enough numbers. This may look in principle as a subitizing-like behavior for small numbers. However, SEB does not predict subitizing in principle. Subitizing-like behavior in SEB results from smaller estimation biases and variability by the observer which may be experimentally imperceptible after rounding to generate a discretized response. (iv) For numbers beyond the subitizing-like range, based on the properties of the log-normal distribution, it can be shown that the coefficient of variation (S1 Note)

$$\frac{SD[\hat{n}]}{E[\hat{n}]} = \sqrt{e^{\hat{\sigma}^2(t)} - 1} \quad (14)$$

does not depend on the input numerosity  $n$ , thus delivering the property of scalar variability, irrespective of  $n$  [5], but here we show that this coefficient will depend on time exposure  $t$ , with the predicted constant coefficient of variation decreasing as  $t$  gets larger proportionally with  $\sqrt{e^{\hat{\sigma}^2(t)} - 1}$  (Fig 1c).

### A thermodynamically inspired model of bounded rationality

Here we briefly introduce a popular approach to studying systems with bounded capacity across domains in human cognition and machine learning: a thermodynamically inspired formalization where information processing is modeled as changes from a default state, which come at some energetic cost, that can be quantified by differences in free energy. This class of models can be applied for the case where an observer intends to minimize some form of expected loss (here we study the case of estimation error minimization), subject to information constraints [17]. More formally, let  $q(\hat{n})$  be a default state (distribution) over

possible responses  $\hat{n}$  in a given environment or context. When presented with a stimulus  $n$ , the resource-constrained observer attempts to transform the initial state  $q$  into a new state of possible responses  $p(\hat{n} | n)$ . This transformation of states can be modeled as the optimization of the free energy functional

$$F[p(\hat{n} | n)] := - \underbrace{\mathbb{E}[L; \hat{n}]}_{\text{Expect. Loss}} - \frac{1}{\beta} \underbrace{D_{KL}(p(\hat{n} | n) \| q(\hat{n}))}_{\text{Constrained State Change}} \quad (15)$$

where  $L$  is a loss function, for instance, the squared error  $(\hat{n} - n)^2$ . The second term is the Kullback-Leibler divergence between  $q$  and  $p(\hat{n} | n)$ , where  $\beta$  trades off the relative importance of changing from the default state  $q$ , thus determining the resources that the observer invests in the estimation task. The goal is to find the optimal distribution of responses

$$p^*(\hat{n} | n) := \arg \max_{p(\hat{n}|n)} F[p(\hat{n} | n)] . \quad (16)$$

The optimal distribution of responses in this variational problem has an analytical solution of the form

$$p^*(\hat{n} | n) \propto q(\hat{n}) \exp(-\beta h(L_n(\hat{n}))) , \quad (17)$$

where  $h$  is a function of  $L$  and potentially other elements incorporated in the expected loss function in Eq 15.

**TIM applied to numerosity estimation.** A recent work applied a model from the TIM family to study a resource-constrained model of human numerosity estimation [21]. This is also a formulation of how the distribution of reported numerosity estimates  $\hat{n}$  of a stimulus magnitude should vary depending on the true stimulus  $n$ . This can be stated generally as the hypothesis that conditional on  $n$  the response distribution  $p(\hat{n} | n)$  is the probability distribution over a set of possible responses  $N$  that minimizes the mean squared error (MSE), subject to the constraint

$$D_{KL}(p_n \| q) \leq C(t) = \min(Rt, B) , \quad (18)$$

where  $C(t)$  is a positive bound that depends on the amount of time  $t$  for which the stimulus is presented. This formulation can be interpreted as a model in which errors in the observer's responses can be attributed to a "cost of control" of the responses: it is difficult for the observer to give responses different from the default state  $q$ , though their response distribution to the individual stimulus  $n$  is optimal given a constraint on the possible precision of their responses.

Similar to our SEB model, in TIM it is assumed that perception extracts information linearly in time at a rate  $R$  until an overall capacity bound  $B$  is reached. The goal is to find the distribution of numerosity estimates  $p^*(\hat{n} | n)$  that minimizes the mean squared error

$$\text{MSE} \equiv \sum_n q(n) \sum_{\hat{n}} p(\hat{n} | n) (\hat{n} - n)^2 \quad (19)$$

under the constraint given in Eq 18.



The optimization problem described above yields the following analytical solution [21]

$$p^*(\hat{n} | n) \propto q(\hat{n}) \exp(-\beta_n q(n)(n - \hat{n})^2), \quad (20)$$

where  $\beta_n$  is chosen to satisfy the bound in Eq 18. Note that this solution has the familiar form obtained in Eq 17 with  $L$  as the loss function  $L_n = (n - \hat{n})^2$ .

While this solution is usually linked to a bounded-rational Bayesian computation (given the observation that the default distribution  $q$  is multiplied by a function of  $n$  given  $\hat{n}$ ), here we clarify that this solution does not correspond to a Bayesian inference process with noisy sensory percepts. In fact, the TIM formulation assumes that the perception of the sensory stimulus  $n$  is noiseless, and all the variability observed during the estimation process is related to the cost of acting accurately, that is, a cost in the precision of response selection when shifting away from the default state. Note that this is fundamentally different from the SEB model, in which all the estimation variability is attributed to noisy sensory encoding.

### Overview of the constraint parameters of the SEB and TIM models

One of our goals is to formalize and make transparent the different elements that play a role in a noisy information transmission process under our model specification, namely, (i) time, (ii) precision of instantaneous information processing, and (iii) energy required to transmit decodable information.

On the one hand, the constraint  $B$  is similar in spirit to the DDMS applied to cognitive processes, where the decision maker is willing to invest a maximum amount of time in processing information due to opportunity costs, in principle irrespective of how precise the instantaneous information processing is. Thus, formally, the time invested in acquiring information is given by

$$\tau = \min(t, t_{max}), \quad (21)$$

with  $t_{max}$  defining the maximum information bound ( $B = R\sqrt{t_{max}/\sigma}$ ).

On the other hand, the parameter  $\Omega$  in our work imposes a constraint on the cost of a decodable transmitted message considering deviations from a status quo state, where energy needs to be injected to transmit decodable information (which in our model specification is given by  $m(n)$ , for any given  $\tau$ ). However, the message  $m$  is not noiseless, where instantaneous processing noise is given by  $\omega$  (otherwise one could make the length of  $m$  infinitesimally small). This means that there exists a natural trade-off between the fidelity of information processing, how much energy the decision maker is willing to (or can) “pay” to disentangle information, and also how much time should be invested in this process. In addition to this, there is an objective that we assume the decision-maker would like to optimize for: minimize the mean squared error. Here it is important to emphasize two points: (i) specifying the model in this way makes transparent the different limitations that the decision-maker must trade-off; (ii) the constraint  $B$  is not a necessary requirement to obtain the optimal solution. It is set following the common knowledge that processing time leads to opportunity costs, in principle irrespective of information processing precision. These constraints will almost surely trade-off in the presence of imprecise information processing (as is also the case in  $n$ -alternative choice DDMS). Should time be allowed to be infinite then decoding would be nearly perfect even with imprecise instantaneous information processing (Eqs 7 and 10); analogously in  $n$ -alternative DDMS with noisy drifts, bounds would be infinitely large, thus converging to errorless decisions.

These constraint assumptions are also present in the TIM:

$$D_{KL}(p_n \| q) \leq C(t) = \min(Rt, B). \quad (22)$$

Here the constraint is similarly given by two parameters: a linear information “processing” rate  $R$ , and (ii) a bounded rate of information processing  $B$ . However, as we discuss in detail in our work, the TIM specification does not consider noisy encoding.

### General similarities and differences between SEB and TIM

We elaborated an illustrative example that allows the predictions of the two models to be solved analytically, thus allowing us to understand the key differences between them (S5 Note). These analyses reveal some similarities between the predictions of the two models, however, there are also notable differences. First, while both models predict that biases decrease in general for larger viewing times, this decrease occurs sooner for SEB than TIM. Second, for a given input stimulus  $n$ , the two models do not imply that  $\text{var}[\hat{n} | n]$  co-varies with the bias in the same way. As the viewing time  $t$  goes to zero, the Bayesian model implies that the variance should fall to zero; TIM implies that this is the case in which estimates should have the highest variance (equal to the variance of the prior distribution).

These analytical insights were studied over all possible responses in the continuous space and do not directly apply to numerosity estimation in the discrete space (S5 Note). Therefore, we conducted numerical analyses to study whether the same signatures emerge in SEB and TIM when the solutions are restricted over the space of positive integers and to give some intuition by visualizing the differences.

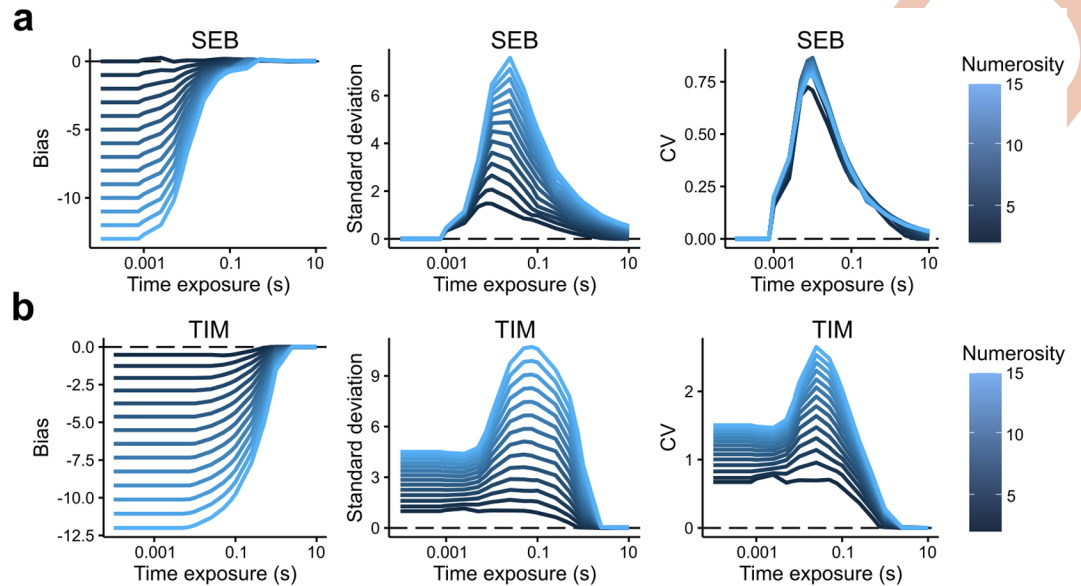
As expected, both models predict that biases decrease in general for larger viewing times, and mirroring the results of the analytical solution, this decrease occurs sooner for SEB than TIM (Fig 2, left panels). However, a fundamental difference between the two models is that as  $t$  goes to zero, the SEB predictions fall to zero, but this is not the case in TIM where the predicted variability of estimates is clearly larger (Fig 2, middle panels).

Finally, the computation of the coefficient of variation ( $\text{CV}[\hat{n}] \equiv \text{SD}[\hat{n}]/E[\hat{n}]$ ) reveals that in SEB this metric is nearly identical for all numerosities  $n$  irrespective of time exposure  $t$ , thus reflecting the scalar variability effect (Fig 2a, right panel). In TIM, however, the scalar variability phenomenon is absent irrespective of time exposure  $t$ . These differences make the two models different and identifiable and generate somewhat different qualitative predictions.

### Efficient numerosity estimation under limited time

We now compare TIM with SEB models using the experimental data of a pre-registered study provided in previous work [21] (see Methods). In brief, on each trial, between 1 and 15 dots were flashed, followed by a noise mask. The participants were then asked to type their estimation of how many dots were displayed. There were three between-participant experiments ( $n=100$  per experiment) that manipulated available stimulus information (variable exposure time:  $t \in [40, 80, 160, 320, 640]$  ms) and different ways of controlling non-numerical properties of the stimuli: the average dot size (experiment 1), the average density of the dots (experiment 2) or the total surface area covered by the dots (experiment 3).

To fully constrain inference solely to the normative solutions of stimulus exposure derived above for both SEB and TIM, we fixed the prior distribution before fitting the behavioral data to a prior equivalent of the form  $1/n^\alpha$  power-law. It has previously been argued that the prior probability of how often numerosities are encountered and represented roughly follows a  $1/n^{\alpha=2}$  power-law distribution [24,25]. Thus, a priori, we choose  $\alpha = 2$ , following the same



**Fig 2. General similarities and differences between SEB and TIM.** **a)** Computation of the bias ( $E[\hat{n} | n] - n$ , left panel), standard deviation ( $SD[\hat{n} | n]$ , middle panel), and the coefficient of variation ( $CV = SD[\hat{n} | n]/E[\hat{n} | n]$ , right panel) as a function of different time exposures  $t$  for different numerosities  $n$  (color scale of the solid lines) in the SEB model. Although we proved that the CV in SEB does not depend on numerosity (14), notice that in the simulations the CV does vary slightly. These slight differences are the result of rounding the expected value of the posterior to the closest integer for the discrete SEB model. For smaller integers, the CV will become more affected by these rounding errors, i.e., due to a slight overestimation of the Bayesian decoder for smaller numbers, the expectation will slightly increase thus dominating the CV ratio. **b)** Same as panel a, but this time computed for TIM. Differences between the two models are particularly salient in the computation of the SD and the CV.

<https://doi.org/10.1371/journal.pcbi.1012790.g002>

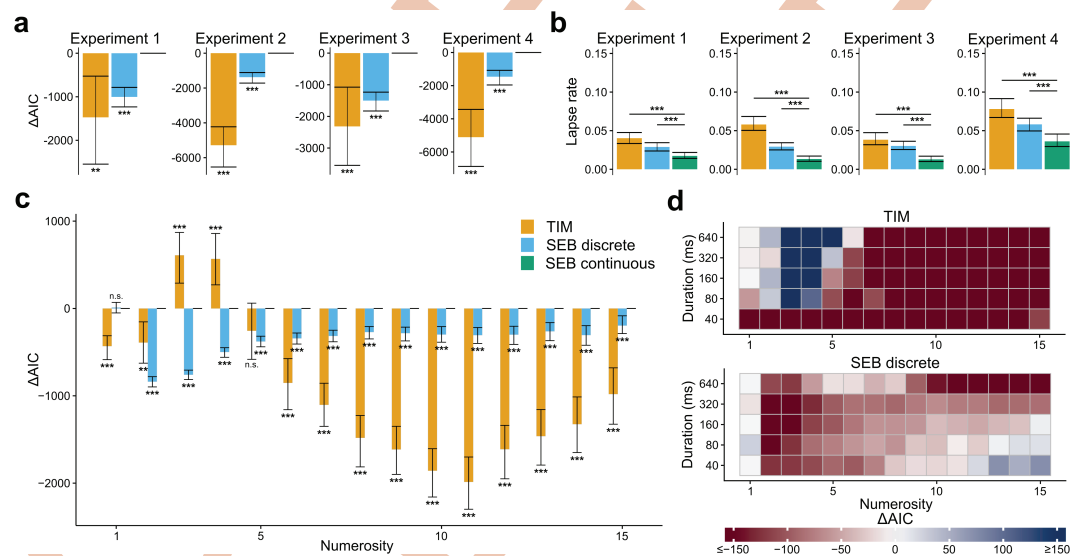
assumption adopted in previous work [21]. By fixing such ecologically valid prior, we alleviate the critique of allowing an arbitrary choice of prior and likelihood functions to fit inference models to the data, as a consequence of which it is sometimes argued that their predictions are potentially vacuous [28]. Nevertheless, it is well possible that each individual has learned their own distribution during their lifespan [29,30]. Therefore, we also considered a more flexible class of models where we allowed the parameters of the prior distribution to be free parameters alongside the capacity constraint and capacity bound.

We considered two possible ways of inferring the numerosity estimates based on the SEB approach (methods): (i) using the analytical solutions over the continuous positive real line, and (ii) using discrete encoding and decoding restricted to the positive integer numbers, thus similar in nature to the TIM specification. Finally, we considered a guessing rate  $g$  in the model fits, which assumes that on  $g$  proportion of trials, participants were distracted and had no information about the number of dots in the display, meaning that their estimate was effectively a random sample from their prior. Thus, both numerosity estimation models SEB and TIM have exactly the same degrees of freedom (the capacity constraint, the capacity bound, and the guessing rate), in addition to the prior parameters in the flexible class of models.

**Quantitative model comparison.** For each experiment where stimulus presentation time  $t$  was manipulated, we fit both types of model to the data of each participant (Methods). In parameter recovery exercises we found that all model parameters are identifiable and this is also confirmed by the weak relationship between parameters across participants (S1 Fig). We

first examined the restricted models where the prior is fixed  $1/n^2$ . For experiment 1 (dot size controlled), we found that the difference in Akaike information criterion (AIC) favoured SEB, where the continuous version of SEB had a clear advantage over TIM:  $\Delta AIC = 1472$  [95%-CI 570-2553] in favor of SEB (paired t-test:  $T(99) = 2.86, p < 0.01, d = 0.29$ ). For experiment 2 (dot density controlled), the difference in AIC is 5284 [95%-CI 4185-6690] in favor of SEB ( $T(99) = 8.70, p < 0.001, d = 0.87$ ). For experiment 3 (dot area controlled), the difference in AIC is 2316 [95%-CI 1218-3686] in favor of SEB ( $T(99) = 3.61, p < 0.001, d = 0.36$ , see Fig 3a). In addition, the SEB continuous model provided better fits than its discrete version ( $T(99) \geq 8.64, p < 0.001, d > 0.86 \Delta AIC \geq 997$ ).

Previous theoretical and empirical work suggests that two ways in which the amount of resources available to process information can be studied are by manipulating time exposure and also by changing stimulus contrast [13]. Thus, we also considered this alternative way of manipulating sensory reliability, which should affect the channel capacity transmission (see Eq 10). To test this, we analyzed data of a numerosity estimation experiment, where in each trial the visual contrast of numerosity was manipulated at a constant presentation time ( $n=100$  participants, experiment 4, Methods). We found that also in this experiment the SEB continuous model fits the data better than TIM ( $\Delta AIC = 5106$ ; [95%-CI 3452-6880]



**Fig 3. The SEB model quantitatively outperforms the TIM model when the prior parameters are fixed.** **a)** Difference in AIC between the SEB continuous model (green) and the TIM model (orange) or the SEB discrete (blue). The  $\Delta AIC$ s were computed for each participant and summed. The error bars represent the 95% confidence interval based on bootstrapping of the participants'  $\Delta AIC$ s. The SEB continuous model outperforms both the TIM model and the SEB discrete model. **b)** Average guessing rate parameter per participant for each model and experiment. The error bars represent the 95% confidence interval based on bootstrapping of the participants' guessing rate. The guessing rate of the SEB continuous model is lower than the TIM model and the SEB discrete model. These results indicate that less variability is associated to lapses of attention in the SEB continuous model, which suggests a better fit to behavior. **c)** Difference in AIC between the SEB continuous model and the TIM model or the SEB discrete model for each numerosity. The error bars represent the 95% confidence interval based on bootstrapping of the participants' AICs. The SEB continuous model outperforms the TIM model except for numerosities 3, 4 and 5 and the SEB discrete model for all numerosities except numerosity 1. **d)** AIC differences between the SEB continuous model and the TIM model (top) and the SEB discrete model (bottom) for all experiments shown for different numerosities and levels of sensory evidence (stimulus presentation duration or contrast). Duration values are assigned to Weber contrasts of experiment 4 for pooling purposes (40ms–10%, 80ms–20%, 160ms–40%, 320ms–80%, 640ms–160%). The SEB continuous model outperforms the TIM and SEB discrete models for most numerosities and levels of sensory evidence.

<https://doi.org/10.1371/journal.pcbi.1012790.g003>

( $T(99) = 5.79, p < 0.001$ ),  $d = 0.58$ , Fig 3a) and the discrete version of SEB ( $\Delta AIC = 1453$ ; [95%-CI 1059-1907] ( $T(99) = 6.69, p < 0.001, d = 0.67$ )).

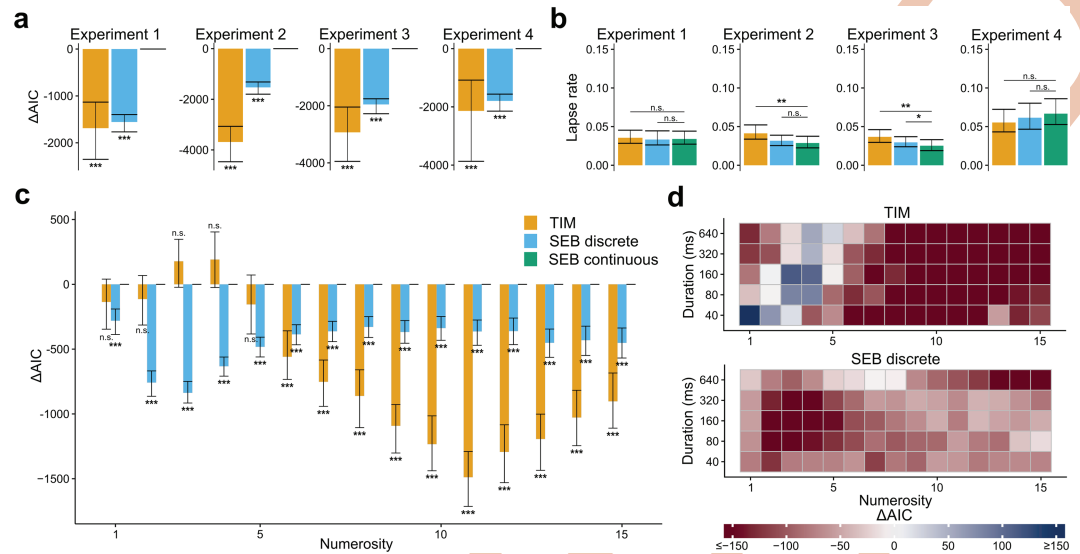
To make sure that the overall quantitative differences were not driven by a few numerosities, we computed the difference in AIC for each numerosity and each model. We found a significant interaction models\*numerosity of the  $\Delta AICs$  ( $F(28, 16758) = 7.84, p < 0.001$ ) with posthoc tests revealing that this effect was more pronounced for higher numerosities (SEB continuous vs. TIM: paired t-tests  $p < 0.001$  for numerosities  $n > 5$ , Fig 3c) and also for  $n \in [1, 2]$  (paired t-tests  $p < 0.01$ ). The relative advantage of the TIM model for  $n \in [3, 4]$  at large presentation times  $t$  might be explained by the fact that smaller numerosities are close to the subitizing range and therefore most of the posterior density mass is concentrated around the input  $n$ , which is better explained by the TIM model as this model has a tendency to subitize more strongly at small numerosities [21]. Interestingly, for  $n \in [1, 2]$ , the Bayesian model predicts noisier estimations (in particular for smaller exposure times  $t$ ) which are not supported by the TIM, with the AICs favoring the former.

Additionally, we inspected the AIC differences split by both numerosity and sensory evidence (time or contrast), finding a similar pattern, but the differences were larger for small levels of sensory evidence (Fig 3d). Thus, SEB appears to be more sensitive to capturing behavior for stimuli generating higher noise levels in the encoding operations.

Moreover, we compared the guessing rates  $g$  between the two kinds of models. Guessing rates can capture unassigned variance in misspecified models, thus we conjectured that a relatively smaller value of  $g$  would provide further evidence for better mechanistic fits captured by the best model. While the guessing rates are overall small (suggesting that the amount of distractions during task performance was minimal), we found that guessing rates were systematically smaller in the SEB model ( $T(99) \geq 5.75, p < 0.001, d > 0.58$  for each experiment, Fig 3b and S1 Table). Thus, while the effects of distraction are estimated to be relatively small in both models, our analyses provide a clear indication that potentially unassigned variance due to distraction is lower in the SEB model relative to TIM. We also note that the information bound parameter  $B$  can be fit for both models which is empirical evidence that this information bound is present in this experimental setup, even for relatively short presentation times (640ms).

We repeated the same set of analyses treating parameters of the log-normal prior as free parameters. The results of these analyses mirrored the initial analyses. That is, (i) we found that the SEB model fit the data better than TIM in all four experiments ( $T(99) \geq 3.54, p < 0.001$ ), Fig 4a), (ii) the continuous version of SEB performed better in general than its discretized version (Fig 4a) and (iii) the guessing rates were significantly smaller in the SEB model than the TIM model for experiments 2 and 3 ( $T(99) \geq 2.95, p < 0.01$  but not for experiments 1 and 4, Fig 4b).

The next question to ask is whether the models with free prior parameters outperformed the models with the prior fixed to  $1/n^2$ . We found that for each model considered here, the models with free prior parameters outperformed their corresponding version with fixed parameters ( $T(99) \geq 5.62, p < 0.001$ , S2 Table). Additionally, to account for population variability in the quantitative metrics between participants across all models considered here, we applied a Bayesian Model Selection which revealed that the Bayesian model with free prior parameters is clearly favored relative to all the other models for experiments 1, 2 and 3 ( $P_{xp} > 0.99$  for each experiment) but equally favored to the TIM model with free prior parameters for experiment 4 ( $P_{xp} = 0.50$ ). These results allow us to conclude two important points. First, variability in the prior parameters of the prior distribution is key to more accurately explaining human numerosity estimations. Second, our results provide a clear indication that the



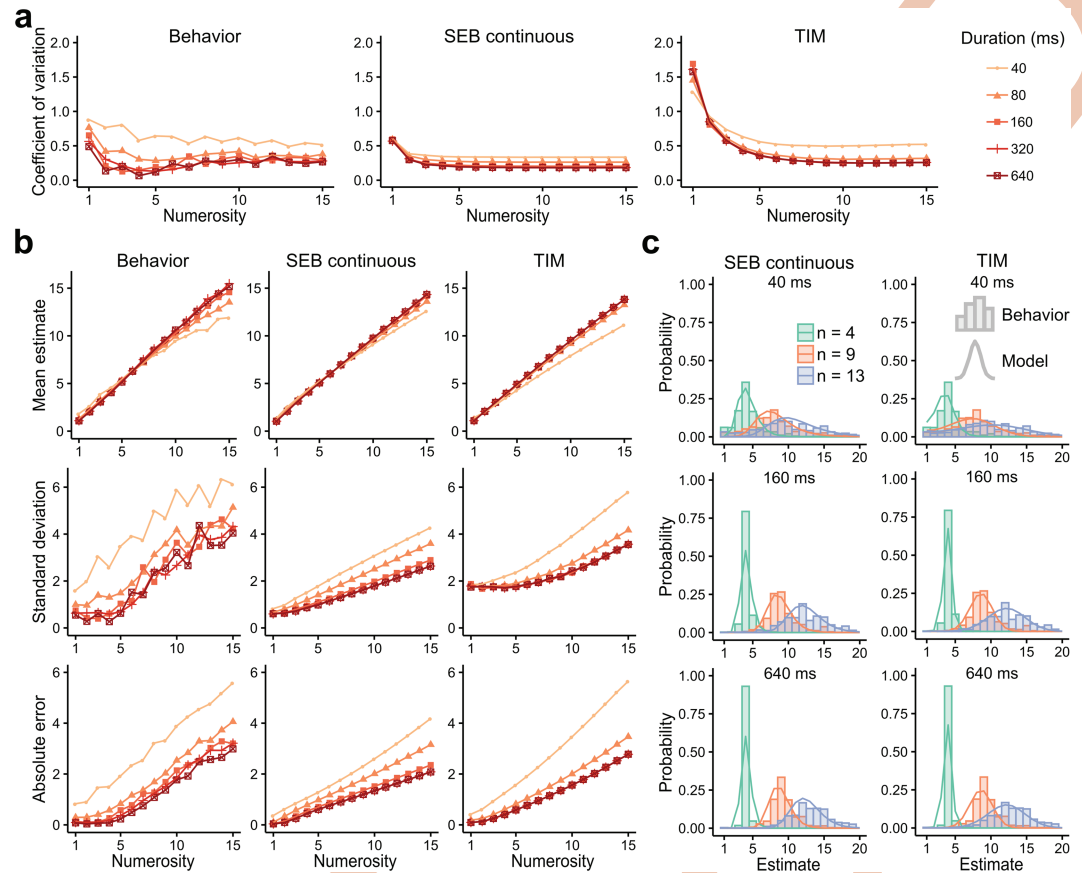
**Fig 4. The SEB model quantitatively outperforms the TIM model when the prior parameters are free.** **a)** Difference in AIC between the SEB continuous model (green) and the TIM model (orange) or the SEB discrete (blue). The  $\Delta AIC$ s were computed for each participant and summed. Error bars represent the 95% confidence interval based on bootstrapping of the participants'  $\Delta AIC$ s. The SEB continuous model outperforms both the TIM model and the SEB discrete model. **b)** Average guessing rate parameter per participant for each model and experiment. Error bars represent the 95% confidence interval based on bootstrapping of the participants' guessing rate. The guessing rate of the SEB continuous model is lower than the TIM model for experiments 2 and 3 but not for experiments 1 and 4 and the SEB discrete model for experiment 3 but not for the other experiments. **c)** Difference in AIC between the SEB continuous model and the TIM model or the SEB discrete model for each numerosity. Error bars represent the 95% confidence interval based on bootstrapping of the participants' AICs. The SEB continuous model outperforms the TIM model except for numerosities 1 to 5 and the SEB discrete model for all numerosities. **d)** AIC differences between the SEB continuous model and the TIM model (top) and the SEB discrete model (bottom) for all experiments shown for different numerosities and levels of sensory evidence (stimulus presentation duration or contrast). Duration values are assigned to Weber contrasts of experiment 4 for pooling purposes (40ms–10%, 80ms–20%, 160ms–40%, 320ms–80%, 640ms–160%). The SEB continuous model outperforms the TIM and SEB discrete models for most numerosities and levels of sensory evidence.

<https://doi.org/10.1371/journal.pcbi.1012790.g004>

effects of temporal time exposure are better captured by the noisy encoding model (SEB) relative to an action control-like model (TIM).

**Qualitative predictions.** We first examined the qualitative features of scalar variability in both data and the predictions of the SEB continuous and the TIM models with free prior parameters. For each numerosity value, we computed the coefficient of variation (CV:  $SD[\hat{n}]/E[\hat{n}]$ ). We found that the empirical data follows the previously observed properties of scalar variability for numerosities greater than 4 (i.e., a flat CV irrespective of numerosity and sensory evidence), with a slight systematic increase of CV for smaller numbers (Fig 5a left). This relative CV increase for small numbers could be explained by the presence of small guessing rates  $g$  which have a greater impact on the CV for small  $n$ . We found that the SEB model accounts for these qualitative observations (Fig 5a middle), however, the TIM model generates slightly different predictions (Fig 5a right).

We found that patterns of estimation biases and variability during numerosity estimation as a function of sensory evidence were in general more closely captured by the SEB relative to the TIM model (Fig 5b top and middle panels). As predicted by our analytical analyses (S5 Note) the rate of increase in noise as a function of  $n$  is larger for the TIM model relative to the SEB model, with the empirical data more closely agreeing with the SEB model. Additionally, given that the TIM model generally requires larger values of guessing rates  $g$  to explain



**Fig 5. The SEB continuous model with free prior parameters qualitatively explains behavior.** **a)** Coefficient of variation ( $SD[\hat{n}]/E[\hat{n}]$ ) of the behavior data (left) and predictions of the SEB model (middle) and TIM model (right) using a prior with free parameters for different numerosities and stimulus presentation duration. Predictions were performed by taking for each parameter the value with this highest density across participants. Duration values are assigned to Weber contrasts of experiment 4 for pooling purposes (40ms–10%, 80ms–20%, 160ms–40%, 320ms–80%, 640ms–160%). The TIM model predicts a higher CV for lower numerosities. This feature is not present in the behavior data nor the SEB predictions. **b)** Mean estimate (top), standard deviation (middle) and absolute error (bottom) of the behavior data (left) and predictions of the SEB model (middle) and TIM model (right). **c)** Posterior distribution of estimates to numerosities 4 (green), 9 (red) and 13 (blue) of the SEB (left) and the TIM (right) model for different stimulation presentation durations (40ms (top), 160ms (middle), 640ms (bottom)). Behavior of participants is shown as histograms.

<https://doi.org/10.1371/journal.pcbi.1012790.g005>

variance, for small  $n$  it predicts larger SDs relative to SEB and empirical data (with a similar pattern for the case of the CV, Fig 5a). However, there is an exception at 40ms, where TIM captures better the range of the standard deviation (from around 2 to 6). Another point where the TIM model appears to do a better job relative to the SEB model is for the absolute error estimations (Fig 5b bottom). Subitizing is more pronounced for low numbers in general, and this reduces both biases and errors for  $n < 5$ . However, beyond the subitizing range and for levels of noise that challenge sensory perception, the SEB model does a better job at capturing all descriptive statistics. To visualize the nature of these differences, the posterior distribution of estimates for both models are shown in Fig 5c for different numerosities and presentation times.

## Discussion

Our theoretical and empirical tests provide clear evidence that a model of Bayesian decoding of noisy internal representations—which provides a normative explanation for the property of scalar variability and can be parsimoniously connected to a theory of limited informational capacity—provides a better account of numerosity estimation data in humans relative to the alternative TIM model considered here. We emphasize that both models: (i) are optimized for the same assumed objective (minimizing the MSE of the estimates), (ii) can be compared under the same assumption about the prior distribution, and (iii) have identical degrees of freedom. Thus, qualitative and quantitative differences between the two information-theoretical models cannot be explained by differences in model complexity, but instead reflect differences in the mechanistic assumptions of the numerosity processing operations. In particular, it is important to note that assumptions about potential encoding and decoding operations are explicitly stated in the Bayesian model. In contrast, these remain “hidden” in the alternative TIM model.

One of our main goals in the development of our modeling framework was to develop an encoding-decoding model incorporating various aspects of human cognition with many antecedents in the literature, which include Brownian motion during evidence processing over time [15] and logarithmic internal representation of numerical quantities [16]. While our proposed model accounts for key qualitative features of the human behavioral data with minimal degrees of freedom, we do not claim that the log-encoding model necessarily accounts for all aspects of numerosity estimation behavior. Indeed, the encoding and decoding strategies that humans and other animals use need not be the same in all contexts [31]. It is equally possible that numerosity processing mechanisms depend on the task at hand, and draw upon an ensemble of strategies that optimize performance under different situations [32,33]. For instance, in future work, it will be interesting to investigate whether situations that involve explicit numerosity estimation vs. discrimination rely on similar or distinct encoding strategies and inference processes.

We assumed that participants employ a log-normal (or power-law) prior, however, it is important to note that the numerosities presented to the participants were drawn from a uniform distribution. We thus implicitly assumed that participants did not rapidly adapt their encoding operations, which might be a reasonable assumption given that participants were not exposed to the new prior for an extended period of time. In addition, in one version of the model fits we allowed the parameters of the prior to be adjusted, resulting in non-uniform distributions, which at the very least suggests that participants did not fully adapt to a uniform prior.

A natural consequence of our theory is that the SEB model parsimoniously endogenizes parameters of the prior distribution in its encoding operations. A testable prediction is that larger prior distribution ranges should lead to more noisy estimates and therefore poorer discriminability for a given capacity bound. This prediction is confirmed by a recent study where it is shown that human participants adapt their numerosity sensitivity for different numerosity ranges, with important implications for risk behavior [34]. Thus two of the key predictions of our theory hold: for a fixed capacity bound sensory reliability should change as a function of (i) time exposure to the sensory stimulus as shown in this study, and (ii) the range of the prior distribution [34].

Additionally, our model predicts that the crossover point from overestimation to underestimation should change as a function of the numerosity range. In this work, we only present data with a fixed range of 1 to 15, thus we cannot test this prediction. However, a previous study using larger numerosity ranges (e.g., up to 30 or 100) found that the cross-over point is



larger for wider numerosity ranges, and crucially, the degree of over- and under-estimation depended on the attentional resources dedicated to numerosity estimation [4]. This result is again in line with the qualitative predictions of our model. Future research could test these predictions quantitatively.

The general working framework of SEB is a Bayesian decoder where the likelihood function optimally endogenizes information in the prior distribution for a given capacity constraint and stimulus exposure time. Here it is important to emphasize that this framework is not restricted to a specific form of the prior. Also, the formulation of processing time, the information processing constraint, and the objective are specified in a general form. However, a fundamental aspect of the SEB model is the specific form of the drift-diffusion term which employs an affine (log-)linear transformation of the input stimulus. The question here is whether this affine function is valid and generalizes for the case of numerosity estimation when it is assumed that the prior distribution changes (i.e., ceases to be power-law or log-normal). We argue that unless there is extensive training over long periods (e.g., many days of ecologically valid adaptation in the absence of any power-law or log-normal distributions), the encoding function may remain of the family of an affine (log-)linear or power-law transformation. While we acknowledge that this argument needs to be tested empirically, recent work appears to support this notion: Prat-Carrabin and Gershman [35] show that when either the prior or objective functions (or incentives) are manipulated during a numerosity estimation task, and even if behavior shows signatures of such adaptation, modeling analyses suggest that subjects' responses feature in all cases logarithmic encoding while the Bayesian decoder takes into account the prior distribution and objective function (incentives). Nevertheless, this form of the encoder might not be optimal and may need to be adapted to the prior distribution and objectives of the organism to achieve efficiency. This idea was studied in previous work both theoretically and empirically where it is shown that if organisms can adapt their encoding functions, they must do so at the earliest stages of sensory processing [36,37], otherwise, information that is lost in the early sensory processing streams cannot be recovered via downstream operations [36,38].

In addition, given the relatively short stimulus presentation times, we assumed that the participants gather information about the stimulus until they reach the time bound  $t_{max}$  or the stimulus disappears. Other specifications of perceptual decision-making problems include endogenous stopping, in which case the observers decide by themselves when to stop gathering information and respond. Although this question is out of the scope of this study, future researchers could build on our model and specify a utility function which takes into account both the rewards for accurate answers and costs related to the time of the decision.

Our model is agnostic about the biological meaning of its parameters. Future research could try to relate them to neural processes [39–42]. We speculate that the bound on the drift rate  $\Omega$  could be related to the information capacity of sensory or evidence accumulation areas (for example how many neurons are used to represent a stimuli or how much precision can these neurons use [43]). This contrasts with the information bound  $B$ , related to the maximum amount of information that can be represented, which could be related to neurons resources in higher cortical areas such as the dlPFC as well as premotor areas [43,44].

Taken together, our findings suggest the fruitfulness of studying optimal models with resource limitations, which can serve as a departing point to understand the neuro-computational mechanisms underlying human behaviour without ignoring the fact that biological systems are limited in their capacity to process information [36,37,45,46]. This highlights that understanding behavior in terms of its objectives while taking into account cognitive limitations, alongside encoding, decoding, and inference processes is likely to be essential to elucidate the mechanisms underlying human cognition.

## Materials and methods

### Participants, data, and experiments

In this work we re-analyzed the data of experiments collected in previous work [21]. In brief, on each trial, between 1 and 15 dots were flashed, followed by a noise mask. The participants were then asked to type their guess of how many dots were displayed. The participants were recruited and carried out the experiment online. There were three between-participant experiments ( $n = 110$  per experiment) that manipulated available stimulus information (variable exposure time:  $t \in [40, 80, 160, 320, 640]$  ms) and different ways of controlling non-numerical properties of the stimuli: the average dot size (experiment 1), the average density of the dots (experiment 2) or the total surface area covered by the dots (experiment 3).

We also studied a fourth experiment ( $n = 110$ ) in which time exposure  $t$  was fixed across trials, but instead display contrast of the dot arrays was varied from trial to trial (experiment 4). In this experiment, the colors of the dots varied between the background (grey) and pitch black, by Weber contrasts of 10%, 20%, 40%, 80% and 160%, at a constant presentation time of  $t = 200$  ms.

Each participant was presented with each combination of numerosity and sensory evidence twice for a total of 150 trials per participant.

### Models

Here we fit the two families of models described in the main text to the data of each participant: (i) We fit the SEB model assuming a log-normal prior with power parameter  $\alpha = 2$ . We fit a continuous version of the model based on the analytical solutions derived in [S1 Note](#) and [S2 Note](#), and a discrete version of this model based on numerical simulations. (ii) Following the procedures of previous work [21], we fit the TIM model assuming a power-law prior with power parameter  $\alpha = 2$ . For both families of models, we also fit a version where the parameters of the log-normal prior were allowed to be free parameters. We also note that analytical solutions in SEB were derived in the continuous space due to mathematical tractability ([S1 Note](#), [S2 Note](#), and [S3 Note](#)). Thus, in order to define the likelihood function of this model in the integer space, we normalized the log probability of estimators (Eq 31) in the integer range  $n \in [1, 2, 3, \dots, 100]$ . Note that both SEB and TIM have exactly the same degrees of freedom ( $R$ ,  $B$ , and  $g$ ), where  $g$  is a guessing rate based on the probability of randomly drawing a value from the default distribution. AQ1

### Quantitative and qualitative analyses

Participants who completed less than 90% of the trials were excluded. Similar to previous work [21] we selected the 100 best participants for each experiment. In addition, trials in which the participant's response was 10 times higher than the presented numerosity or the response time was superior to 10s were excluded. This additional data cleaning leads to the rejection of 142 trials out of 14,997 for experiment 1, 143 out of 14,993 for experiment 2, 172 out of 15,000 for experiment 3 and 187 out of 15,000 for experiment 4. Each model was fit individually to each participant using the DEoptim package [47] in the statistical language R [48] with a number of iterations set to 100. The limits for the parameter search space were set to (0.1,200) for  $R$ , (0.1, 20) for  $B$  and (0.0001, 0.5) for  $g$ . In the models where the prior was free, the search space of the prior parameters was (-50,50) for  $\mu$  and (0.1,100) for  $\sigma$ . Model comparison was performed based on the Akaike information criterion (AIC). Using other model comparison metrics such as the Bayesian information criterion (BIC) does not change the conclusions of our work.

In Figs 3 and 4 and main text, we report the sum of the AIC difference relative to the best model across participants for each experiment, and report the 95% bootstrap confidence interval (95%-CI). We also computed two-sided paired t-tests based on the AICs obtained for each participant between the SEB and the TIM models. Likewise, we computed two-sided paired t-tests based on the guess rate parameter  $g$  obtained from each participant in the SEB model relative to the guess rates obtained in the TIM model. We report effect sizes as Cohen's  $d$ . The qualitative predictions were computed based on the value with the highest density for each parameter at the population level. Each statistic was computed separately for each experiment and then averaged across experiments.

Details regarding the theoretical derivations of the SEB model and the analytical comparison between TIM and SEB models are given in detail in the Supplementary Notes (S1 Note, S2 Note, S3 Note, S4 Note and S5 Note).

### Ethics statement

All data analysis from human participants is based on an openly available dataset [21], therefore no Institutional Review Board approval was required for this study.

### Supporting information

**S1 Note. General specification and derivation of the logarithmic noisy encoding and Bayesian decoding model.**

(PDF)

**S2 Note. Logarithmic noisy encoding and Bayesian decoding under limited informational capacity and temporal sensory exposure.**

(PDF)

**S3 Note. Recapitulation of the Gaussian channel capacity derivation.**

(PDF)

**S4 Note. Influence of the prior on the crossover point.**

(PDF)

**S5 Note. Comparison between the family TIM and SEB models.**

(PDF)

**S1 Fig. Correlations of parameter fits.**

(PDF)

**S1 Table. Parameter fits.**

(PDF)

**S2 Table. AIC of the model fit.**

(PDF)

### Author contributions

**Conceptualization:** Joseph A. Heng, Michael Woodford, Rafael Polania.

**Data curation:** Joseph A. Heng, Rafael Polania.

**Formal analysis:** Joseph A. Heng, Michael Woodford, Rafael Polania.

**Funding acquisition:** Michael Woodford, Rafael Polania.

**Investigation:** Joseph A. Heng, Michael Woodford, Rafael Polania.

**Methodology:** Joseph A. Heng, Michael Woodford, Rafael Polania.

**Project administration:** Rafael Polania.

**Resources:** Rafael Polania.

**Supervision:** Rafael Polania.

**Validation:** Joseph A. Heng.

**Visualization:** Joseph A. Heng.

**Writing – original draft:** Joseph A. Heng, Michael Woodford, Rafael Polania.

**Writing – review & editing:** Joseph A. Heng, Michael Woodford, Rafael Polania.

## References

1. Nieder A. The adaptive value of numerical competence. *Trends Ecol Evol.* 2020;35(7):605–17. <https://doi.org/10.1016/j.tree.2020.02.009> PMID: 32521244
2. Anobile G, Arrighi R, Castaldi E, Burr DC. A sensorimotor numerosity system. *Trends Cogn Sci.* 2021;25(1):24–36. <https://doi.org/10.1016/j.tics.2020.10.009> PMID: 33221159
3. Revkin SK, Piazza M, Izard V, Cohen L, Dehaene S. Does subitizing reflect numerical estimation? *Psychol Sci.* 2008;19(6):607–14. <https://doi.org/10.1111/j.1467-9280.2008.02130.x> PMID: 18578852
4. Anobile G, Cicchini GM, Burr DC. Linear mapping of numbers onto space requires attention. *Cognition.* 2012;122(3):454–9. <https://doi.org/10.1016/j.cognition.2011.11.006> PMID: 22154543
5. Izard V, Dehaene S. Calibrating the mental number line. *Cognition.* 2008;106(3):1221–47. <https://doi.org/10.1016/j.cognition.2007.06.004> PMID: 17678639
6. Whalen J, Gallistel CR, Gelman R. Nonverbal counting in humans: the psychophysics of number representation. *Psychol Sci.* 1999;10(2):130–7. <https://doi.org/10.1111/1467-9280.00120>
7. Inglis M, Gilmore C. Sampling from the mental number line: how are approximate number system representations formed? *Cognition.* 2013;129(1):63–9. <https://doi.org/10.1016/j.cognition.2013.06.003> PMID: 23831565
8. Jazayeri M, Shadlen MN. Temporal context calibrates interval timing. *Nat Neurosci.* 2010;13(8):1020–6. <https://doi.org/10.1038/nn.2590> PMID: 20581842
9. Navarrete A, van Schaik CP, Isler K. Energetics and the evolution of human brain size. *Nature.* 2011;480(7375):91–3. <https://doi.org/10.1038/nature10629> PMID: 22080949
10. Woodford M. Modeling imprecision in perception, valuation, and choice. *Annu Rev Econ.* 2020;12(1):579–601. <https://doi.org/10.1146/annurev-economics-102819-040518>
11. Bhui R, Lai L, Gershman SJ. Resource-rational decision making. *Curr Opin Behav Sci.* 2021;41:15–21. <https://doi.org/10.1016/j.cobeha.2021.02.015>
12. Louie K, Glimcher PW. Efficient coding and the neural representation of value. *Ann N Y Acad Sci.* 2012;1251:13–32. <https://doi.org/10.1111/j.1749-6632.2012.06496.x> PMID: 22694213
13. Wei X-X, Stocker AA. A Bayesian observer model constrained by efficient coding can explain “anti-Bayesian” percepts. *Nat Neurosci.* 2015;18(10):1509–17. <https://doi.org/10.1038/nn.4105> PMID: 26343249
14. Glimcher PW. Efficiently irrational: deciphering the riddle of human choice. *Trends Cogn Sci.* 2022;26(8):669–87. <https://doi.org/10.1016/j.tics.2022.04.007> PMID: 35643845
15. Ratcliff R. A theory of memory retrieval. *Psychol Rev.* 1978;85(2):59–108. <https://doi.org/10.1037/0033-295x.85.2.59>
16. Nieder A, Miller EK. Coding of cognitive magnitude: compressed scaling of numerical information in the primate prefrontal cortex. *Neuron.* 2003;37(1):149–57. [https://doi.org/10.1016/s0896-6273\(02\)01144-3](https://doi.org/10.1016/s0896-6273(02)01144-3) PMID: 12526780
17. Ortega PA, Braun DA. Thermodynamics as a theory of decision-making with information-processing costs. *Proc R Soc A.* 2013;469(2153):20120683. <https://doi.org/10.1098/rspa.2012.0683>
18. Ortega PA and Stocker AA. Human decision-making under limited time. *Adv Neural Inf Process Syst.* 2016:29.

19. Friston K. The free-energy principle: a rough guide to the brain?. *Trends Cogn Sci.* 2009;13(7):293–301. <https://doi.org/10.1016/j.tics.2009.04.005> PMID: 19559644
20. Wolpert DH. Complex engineering systems, chapter information theory—the bridge connecting bounded rational game theory and statistical physics. Berlin, Heidelberg: Springer; 2006. <https://doi.org/10.1007/3-540-32834-3>
21. Cheyette SJ, Piantadosi ST. A unified account of numerosity perception. *Nat Hum Behav.* 2020;4(12):1265–72. <https://doi.org/10.1038/s41562-020-00946-0> PMID: 32929205
22. Khaw MW, Li Z, Woodford M. Cognitive imprecision and small-stakes risk aversion. *Rev Econ Stud.* 2020;88(4):1979–2013. <https://doi.org/10.1093/restud/rdaa044>
23. LIMPERT E, STAHEL WA, ABBT M. Log-normal distributions across the sciences: keys and clues. *BioScience.* 2001;51(5):341. [https://doi.org/10.1641/0006-3568\(2001\)051\[0341:Indats\]2.0.co;2](https://doi.org/10.1641/0006-3568(2001)051[0341:Indats]2.0.co;2)
24. Piantadosi ST. A rational analysis of the approximate number system. *Psychon Bull Rev.* 2016;23(3):877–86. <https://doi.org/10.3758/s13423-015-0963-8> PMID: 26755188
25. Dehaene S, Mehler J. Cross-linguistic regularities in the frequency of number words. *Cognition.* 1992;43(1):1–29. [https://doi.org/10.1016/0010-0277\(92\)90030-I](https://doi.org/10.1016/0010-0277(92)90030-I) PMID: 1591901
26. Taylor MM, Lindsay PH, Forbes SM. Quantification of shared capacity processing in auditory and visual discrimination. *Acta Psychol (Amst).* 1967;27:223–9. [https://doi.org/10.1016/0001-6918\(67\)99000-2](https://doi.org/10.1016/0001-6918(67)99000-2) PMID: 6062214
27. Cover TM. Elements of information theory. Wiley; 1999. <https://doi.org/10.1002/047174882x>
28. Bowers JS, Davis CJ. Bayesian just-so stories in psychology and neuroscience. *Psychol Bull.* 2012;138(3):389–414. <https://doi.org/10.1037/a0026450> PMID: 22545686
29. Stocker AA, Simoncelli EP. Noise characteristics and prior expectations in human visual speed perception. *Nat Neurosci.* 2006;9(4):578–85. <https://doi.org/10.1038/nn1669> PMID: 16547513
30. Girshick AR, Landy MS, Simoncelli EP. Cardinal rules: visual orientation perception reflects knowledge of environmental statistics. *Nat Neurosci.* 2011;14(7):926–32. <https://doi.org/10.1038/nn.2831> PMID: 21642976
31. Teufel C, Fletcher PC. Forms of prediction in the nervous system. *Nat Rev Neurosci.* 1–12 (2020). ISSN 1471-003X. <https://doi.org/10.1038/s41583-020-0275-5>
32. Heng JA, Woodford M, Polania R. Efficient sampling and noisy decisions. *Elife.* 2020;9:e54962. <https://doi.org/10.7554/eLife.54962> PMID: 32930663
33. Testolin A, McClelland JL. Do estimates of numerosity really adhere to Weber's law? A reexamination of two case studies. *Psychon Bull Rev.* 2020;28(1):158–68. <https://doi.org/10.3758/s13423-020-01801-z>
34. Frydman C, Jin LJ. Efficient coding and risky choice. *Quart J Econ.* 2021;137(1):161–213. <https://doi.org/10.1093/qje/qjab031>
35. Prat-Carrabin A, Gershman SJ. Bayesian estimation yields anti-Weber variability. 2024. <https://doi.org/10.1101/2024.08.08.607196>
36. Schaffner J, Bao SD, Tobler PN, Hare TA, Polania R. Sensory perception relies on fitness-maximizing codes. *Nat Hum Behav.* 2023;7(7):1135–51. <https://doi.org/10.1038/s41562-023-01584-y> PMID: 37106800
37. Grujic N, Brus J, Burdakov D, Polania R. Rational inattention in mice. *Sci Adv.* 2022;8(9):eabj8935. <https://doi.org/10.1126/sciadv.abj8935> PMID: 35245128
38. Polania R, Burdakov D, Hare TA. Rationality, preferences, and emotions with biological constraints: it all starts from our senses. *Trends Cogn Sci.* 2024;28(3):264–77. <https://doi.org/10.1016/j.tics.2024.01.003> PMID: 38341322
39. Polania R, Krajbich I, Grueschow M, Ruff CC. Neural oscillations and synchronization differentially support evidence accumulation in perceptual and value-based decision making. *Neuron.* 2014;82(3):709–20. <https://doi.org/10.1016/j.neuron.2014.03.014> PMID: 24811387
40. Pisauro MA, Fouragnan E, Retzler C, Philiastides MG. Neural correlates of evidence accumulation during value-based decisions revealed via simultaneous EEG-fMRI. *Nat Commun.* 2017;8:15808. <https://doi.org/10.1038/ncomms15808> PMID: 28598432
41. Krajbich I, Mitsumasu A, Polania R, Ruff CC, Fehr E. A causal role for the right frontal eye fields in value comparison. *Elife.* 2021;10:e67477. <https://doi.org/10.7554/eLife.67477> PMID: 34779767
42. DePasquale B, Brody CD, Pillow JW. Neural population dynamics underlying evidence accumulation in multiple rat brain regions. *Elife.* 2024;13:e84955. <https://doi.org/10.7554/eLife.84955> PMID: 39162374
43. Katz LN, Yates JL, Pillow JW, Huk AC. Dissociated functional significance of decision-related activity in the primate dorsal stream. *Nature.* 2016;535(7611):285–8. <https://doi.org/10.1038/nature18617> PMID: 27376476

44. Forstmann BU, Anwander A, Schäfer A, Neumann J, Brown S, Wagenmakers E-J, et al. Cortico-striatal connections predict control over speed and accuracy in perceptual decision making. *Proc Natl Acad Sci U S A*. 2010;107(36):15916–20. <https://doi.org/10.1073/pnas.1004932107> PMID: 20733082
45. Woodford M. Prospect theory as efficient perceptual distortion. *Am Econ Rev*. 2012;102(3):41–6. <https://doi.org/10.1257/aer.102.3.41>
46. Ganguli D, Simoncelli EP. Efficient sensory encoding and Bayesian inference with heterogeneous neural populations. *Neural Comput*. 2014;26(10):2103–34. [https://doi.org/10.1162/NECO\\_a\\_00638](https://doi.org/10.1162/NECO_a_00638) PMID: 25058702
47. Mullen K, Ardia D, Gil D, Windover D, Cline J. DEoptim: AnRPackage for global optimization by differential evolution. *J Stat Soft*. 2011;40(6):1–26. <https://doi.org/10.18637/jss.v040.i06>
48. R Core Team. R: a language and environment for statistical computing. 2019.

# AUTHOR QUERY FORM

AQ1.

Please provide appropriate citation for (Eq. 31) in the sentence starting with “Thus, in order to define the likelihood function...”

UNCORRECTED  
PROOF