

# Bayesian Encoding and Decoding as Distinct Perspectives on Neural Coding

Richard D. Lange, Sabyasachi Shivkumar\*, Ankani Chattoraj\*, Ralf M. Haefner

## Abstract

The Bayesian Brain hypothesis, according to which the brain implements statistically optimal algorithms, is one of the leading theoretical frameworks in neuroscience. There are two distinct underlying philosophies: one in which the brain recovers experimenter-defined structures in the world from sensory neural activity (decoding), and another in which it represents latent quantities in an internal model (encoding). We argue that an implicit disagreement on this point underlies some of the debate surrounding the neural implementation of statistical algorithms, in particular the difference between sampling-based and parametric distributional codes. To demonstrate the complementary nature of the two approaches, we have shown mathematically that encoding by sampling can be equivalently interpreted as decoding task variables in a manner consistent with linear probabilistic population codes (PPCs), a popular decoding approach. Awareness of these differences in perspective helps misunderstandings and false dichotomies, and future research will benefit from an explicit discussion of the relative advantages and disadvantages of either approach to constructing models.

## 1 Introduction

According to the Bayesian Brain hypothesis, one of the main operations of neural circuits is to carry out statistical computations by flexibly combining prior knowledge with new evidence and evaluating quantities of interest with respect to the entire posterior distribution. In the case of perception, prior knowledge is assumed either to come from experience with the world during development or to be encoded genetically having been learned over the course of generations. While any given sensory measurement may be noisy or ambiguous – providing a wide likelihood function in Bayesian terms – prior knowledge is deployed to resolve these ambiguities when possible (von Helmholtz, 1925). The Bayesian framework has been instrumental for our understanding of perception (Knill and Richards, 1996; Pouget et al., 2013).

At the core of the Bayesian Brain hypothesis is the idea that neural activity corresponds to probability distributions rather than point estimates – such schemes are known as “distributional codes” (Zemel et al., 1998). Previous surveys of distributional codes have emphasized a distinction between sampling-based and parametric codes (Fiser et al., 2010; Pouget et al., 2013; Sanborn, 2015; Gershman and Beck, 2016). From a general theoretical standpoint, both sampling and parametric codes have advantages and disadvantages. In the context of neuroscience, sampling and parametric codes have also been compared with respect to the simplicity of implementing computations believed to be important for the brain, such as cue combination and marginalization (Fiser et al., 2010). Further, numerous studies have empirically tested for properties of sampling or parametric codes in neural responses. Sampling codes have been argued to explain spontaneous cortical activity (Berkes et al., 2011), neural variability (Hoyer and Hyvärinen, 2003), structure in noise correlations (Haefner et al., 2016; Bányai et al., 2019), onset transients and oscillations (Aitchison and Lengyel, 2016; Hennequin et al., 2018; Echeveste et al., 2019), and more (Orbán et al., 2016). Meanwhile, parametric codes have been cited in explanations of contrast-invariant tuning (Ma et al., 2006), near-linearity during cue-combination (Fetsch et al., 2011, 2013), evidence integration dynamics in parietal cortex (Beck et al., 2008; Hou et al., 2019), divisive normalization (Beck et al., 2011), and more (Pouget et al., 2013). Importantly, sampling and parametric codes have so far always been discussed and compared as competing

---

\*equal contribution

26 and mutually exclusive mathematical models of the same neural circuits, with no decisive evidence presented  
27 favoring one over the other model.

28 Here, we describe how part of this debate can be resolved by considering that sampling and parametric  
29 codes, as they are usually discussed, reflect two distinct and *complementary* philosophies on how to con-  
30 struct models of inference in the brain. In particular, the primary goal of this paper is to clearly establish  
31 a distinction between what we call **Bayesian Encoding** and **Bayesian Decoding** perspectives on the  
32 Bayesian Brain hypothesis. These two perspectives constitute different ways of thinking about the kinds of  
33 inference problems faced by the brain and over what variables which inference is performed. Not making  
34 these differences explicit has led to confusion about how to interpret neural data. The distinction between  
35 an encoding and a decoding perspective has several components, an understanding of which we hope will  
36 clarify future research.

37 We illustrate the complementary nature of these two philosophies using a toy model, previously presented  
38 at NeurIPS (Shivkumar et al., 2018). In this example, we *construct* a sampling-based encoding over a linear  
39 Gaussian image model (Olshausen and Field, 1996, 1997), and show analytically that firing rates in this  
40 model are equivalent to a Probabilistic Population Code (PPC) over arbitrary scalar stimuli in a task. There  
41 is thus no inherent contradiction in saying that the brain is *both* sampling (in the “Bayesian Encoding”  
42 sense) *and* a parametric code (in the “Bayesian Decoding” sense). We conclude with a discussion of other  
43 possible connections between sampling and parametric codes and distributional neural codes in general.

## 44 2 Results

45 Both Bayesian Encoding and Bayesian Decoding fall under the umbrella of *distributional* neural codes. This  
46 means that any given pattern of neural activity is interpreted not as representing a point estimate of some  
47 quantity, but as representing an entire probability distribution over it. The nature of this “quantity” is key  
48 to the distinction between both frameworks.

### 49 2.1 Bayesian Encoding

50 We define **Bayesian Encoding** as the view that there exists a probability distribution over some quantity  
51 of interest to the brain, and that the primary function of sensory neurons is to compute and represent an  
52 approximation to this distribution. We use the term “encoding” because the probability distribution that is  
53 represented conceptually precedes the actual neural responses. That is, in Bayesian encoding models, there  
54 exists a reference distribution that is defined independently of how neurons actually respond, and which is  
55 approximately encoded by neural responses.

56 The Bayesian Encoding perspective requires a probabilistic model that defines the reference distribution.  
57 In the context of the sensory system, this model often takes the form of an internal generative model of  
58 sensory inputs (Figure 1a). With this perspective, the long-term goal of sensory areas of the brain is to  
59 develop a statistical model of its sensory inputs. Sensory data, such as an image on the retina, are *explained*  
60 as the result of higher order causes. Whereas an image on the retina is high-dimensional and complex, latent  
61 variables tell their story: objects, lights, textures, and optics interacted to create each image. A generative  
62 model makes this process explicit by assigning prior probabilities to the (co)occurrence of latent variables and  
63 by quantifying the likelihood of generating a particular sensory observation from a particular configuration  
64 of latent variables. The encoded distribution in this framework is defined over the variables in this statistical  
65 model.

66 For latent variables  $\mathbf{x}$  and sensory input  $\mathbf{I}$ , optimal inference means computing the posterior distribution,

$$p_b(\mathbf{x}|\mathbf{I}) = \frac{p_b(\mathbf{I}|\mathbf{x})p_b(\mathbf{x})}{p_b(\mathbf{I})}. \quad (1)$$

67 We use the subscript b in  $p_b(\mathbf{x}, \mathbf{I})$  to refer to quantities in the brain’s internal model to distinguish them  
68 from other types of probabilities such as a decoder’s uncertainty. The Bayesian Encoding perspective poses  
69 the question of how neural circuits could compute and represent the posterior distribution  $p_b(\mathbf{x}|\mathbf{I})$  for any  
70 sensory  $\mathbf{I}$ , given the internal model that the brain has learned (Figure 1b). In general, exact inference  
71 is an intractable problem (Murphy, 2012; Wainwright and Jordan, 2008; Bishop, 2006), leading to the

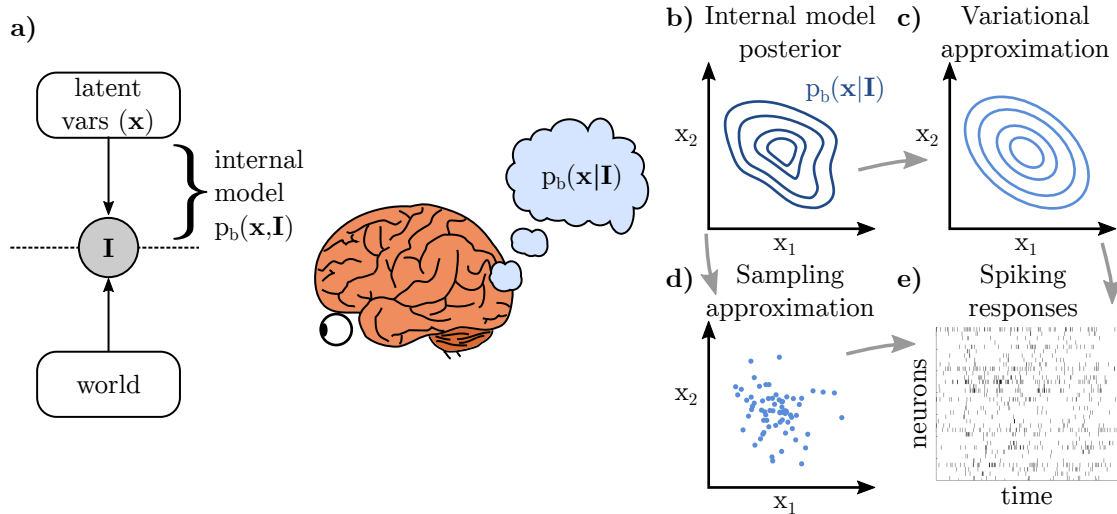


Figure 1: Visualization of Bayesian Encoding. **a)** A common assumption of Bayesian Encoding is that the brain constructs an internal model of the world, and that perceptual inferences are about quantities in the internal model, as opposed to being about external quantities in the world *per se*. This diagram emphasizes this distinction between the world and an internal model. Whether or not stimuli come from natural experience or from an artificial task, the brain computes a posterior over internal variables,  $p_b(\mathbf{x}|\mathbf{I})$ , in all cases. **b-e)** The defining feature of Bayesian Encoding is the existence of a “true” distribution (b), often the posterior over a latent variable,  $\mathbf{x}$ , given a sensory measurement,  $\mathbf{I}$ . One then typically assumes an approximation scheme such as variational inference (b→c) or sampling (b→d), and that this approximation is then realized in patterns of neural activity (e).

72 question of how the brain could compute and represent an *approximation* to the true posterior (Figure 1c-e).  
 73 This line of reasoning motivates work on “neurally plausible approximate inference algorithms,” including  
 74 approaches with connections to sampling-based inference (Figure 1d), as well as approaches inspired by  
 75 variational inference techniques (Figure 1c) (reviewed in Fiser et al. (2010); Sanborn (2015); Gershman and  
 76 Beck (2016)).

## 77 2.2 Bayesian Decoding

78 We define **Bayesian Decoding** as the perspective in which neural activity is treated as *given*, and emphasis  
 79 is placed on the statistical uncertainty of a decoder observing those neural responses. Bayesian Decoding is  
 80 closely related to ideal observer models in psychophysics involving tasks that require the estimation of scalar  
 81 aspects of a presented stimulus (e.g. its orientation or its contrast) or a decision whether the stimulus belongs  
 82 to one of two or more discrete classes (e.g. “left” or “right”). Of course, any stimulus  $s$  that elicits neural  
 83 responses  $\mathbf{r}$  is optimally decoded by computing  $p(s|\mathbf{r})$ . In general, this decoder may be complex or sensitive  
 84 to context or other “nuisance variables.” The key question within the Bayesian Decoding framework is this:  
 85 what conditions must the stimulus-driven neural activity ( $p(\mathbf{r}|s)$ ) fulfill such that the decoder ( $p(s|\mathbf{r})$ ) is both  
 86 simple (e.g. linear) and invariant to changes in context? For instance, linearity and invariance constraints  
 87 on the decoder imply constraints on tuning curves and the distribution of neural noise (Zemel et al., 1998;  
 88 Ma et al., 2006).

89 There is little practical difference between this definition of Bayesian Decoding and familiar notions of  
 90 optimal neural decoding, except in one’s philosophical stance towards inference in the brain, and hence  
 91 in the kinds of problems and tools that are emphasized. Classically, decoding is either a tool for assessing  
 92 information content in neural responses or a mechanistic model of how they impact behavior. In the Bayesian  
 93 setting, one might further invoke the language of ideal observers and priors. However, contrasting Bayesian  
 94 versus classical decoding is not pertinent to our main argument; we are instead interested in the distinction  
 95 of both with Bayesian Encoding.

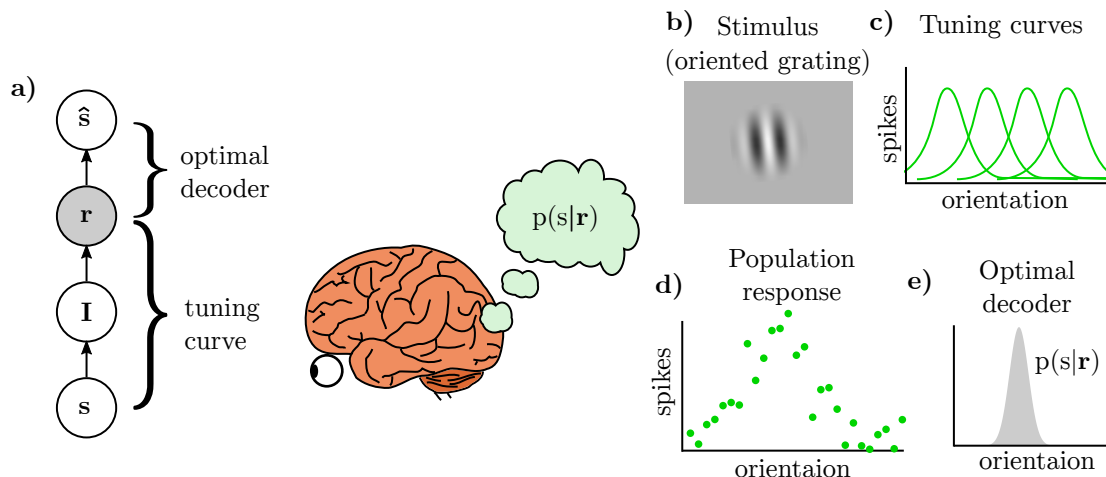


Figure 2: Visualization of Bayesian Decoding. **a)** Decoding is fundamentally a problem of estimating *external* quantities from internal (neural) representations. This diagram emphasizes the symmetry between stimuli that exist in the world, and quantities estimated or inferred in the brain. Here, a scalar stimulus,  $s$ , elicits neural responses,  $r$ , mediated by an image,  $I$ . The decoding question is how the brain forms an internal estimate,  $\hat{s}$ , from  $r$ . **b)** The decoding problem usually begins with a stimulus, such as the direction of motion of dots viewed through an aperture. **c-e)** Given a population of neurons’ tuning curves to  $s$  (c) and an observation of spikes on a single trial (d), an optimal decoder computes  $p(s|r)$  (e). A PPC is a decoder with two convenient properties: it is an exponential family with natural parameters linearly related to  $r$ , and the decoder is invariant to irrelevant nuisance variables if they only scale the tuning curves.

96 Probabilistic Population Codes (PPCs), as introduced by Ma et al (2006), exemplify the Bayesian De-  
 97 coding approach. PPCs imply one way to construct a Bayesian decoder that is both simple and invariant  
 98 to nuisance: if a population of neurons tuned to  $s$  have “Poisson-like” variability, then the optimal decoder  
 99 is part of the exponential family with firing rates as natural parameters. This is a particularly “convenient”  
 100 representation for taking products of two distributions (Ma et al., 2006; Beck et al., 2008). Perhaps even  
 101 more important is the notion of *invariance* afforded by a PPC: as long as nuisance variables such as image  
 102 contrast or dot coherence only multiplicatively scale tuning curves, the decoder can ignore them.

103 Importantly, linearity for cue combination and multiplicative gain by nuisance variables are what con-  
 104 stitute the *predictions* of PPCs. In classical decoding approaches, neural responses are simply “given,” not  
 105 prescribed by a theory. In the Bayesian Decoding framework generally, and in the case of PPCs in particular,  
 106 imposing constraints on the decoder constrain the possible set of evoked response distributions,  $p(r|s)$ . These  
 107 constraints are then formulated as predictions and tested empirically (Fetsch et al., 2011, 2013; Pouget et al.,  
 108 2013; Hou et al., 2019).

## 109 2.3 Contrasting Bayesian Encoding and Bayesian Decoding

110 There are three key differences between the Bayesian Encoding and Bayesian Decoding perspectives involving  
 111 (1) what they assume the brain is inferring, (2) implicit notions of causality, and (3) the empirical data and  
 112 other arguments used to motivate them. As our goal is to summarize and categorize a large and diverse  
 113 sub-field, there will be exceptions to each rule, but we expect these distinctions to be useful for framing  
 114 further discussions.

### 115 2.3.1 Differences in what is assumed to be inferred

116 An integral part of the Bayesian Encoding framework is the existence of an abstract internal model that is  
 117 defined independently of how neurons actually respond. The model is independent of neurons in the sense

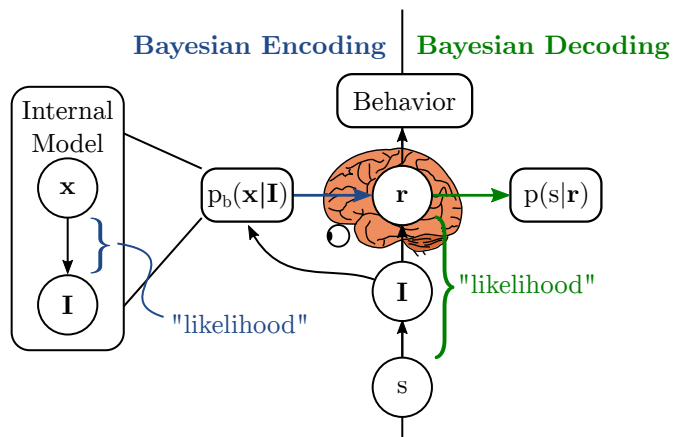


Figure 3: Side-by-side comparison of Bayesian Encoding and Bayesian Decoding. In both frameworks, it is understood that there exists a mechanistic connection between stimuli ( $\mathbf{I}$ ), sensory neural responses ( $\mathbf{r}$ ), and behavior. In the Bayesian Decoding framework, emphasis is placed on the uncertainty of a decoder estimating a (usually scalar) stimulus parameter  $s$  from  $\mathbf{r}$  (green arrow). Bayesian Encoding posits the existence of an internal model with latent variables  $\mathbf{x}$ , and that neural responses ( $\mathbf{r}$ ) encode the computation of a posterior distribution ( $p_b(\mathbf{x}|\mathbf{I})$ ). The blue arrow from  $p_b(\mathbf{x}|\mathbf{I})$  to  $\mathbf{r}$  is an instance of *downward causation*, since changes to the posterior imply changes to neural responses. In Bayesian Decoding, the “likelihood” refers to  $p(\mathbf{r}|s)$ , and the inference problem is to recover  $s$  from  $\mathbf{r}$ . In Bayesian Encoding, the “likelihood” refers to the internal model’s  $p_b(\mathbf{I}|\mathbf{x})$ , and the inference problem is to recover  $\mathbf{x}$  from  $\mathbf{I}$  and to embed the posterior over  $\mathbf{x}$  in  $\mathbf{r}$ .

118 that the same model could in principle be implemented *in silico* or in the brains of other individuals or other  
 119 species. Translating from inference in an internal model into predictions for neural data usually requires  
 120 an additional linking hypothesis on the nature of distributional codes, such as whether neurons sample or  
 121 encode variational parameters, and how either samples or parameters correspond to observable biophysical  
 122 quantities like membrane potentials, spike times or spike counts.

123 The brain’s internal model is typically assumed to have been calibrated through exposure to natural  
 124 stimuli (Berkes et al., 2011) and to only change slowly with extensive exposure to new stimuli. For this  
 125 reason, the generative model in Bayesian Encoding models is often assumed to be task-independent; *what*  
 126 *the brain infers* is assumed to not be under the control of an experimenter. One exception to this rule is a  
 127 family of models in which the *prior* over internal variables changes through extensive exposure to stimuli in  
 128 a particular task (Haefner et al., 2016; Lange and Haefner, 2020).

129 In contrast, the Bayesian Decoding view usually deals directly with estimation of task-relevant variables.  
 130 For instance, in an motion discrimination task, a Bayesian Decoding question would be how the brain  
 131 represents uncertainty over directions of motion. Importantly, answering this question does not require a  
 132 generative model of possible motion stimuli; it requires only a statistical model of the relation between  
 133 scalar motion direction (and possibly nuisance variables like coherence) and neural responses, i.e.  $p(\mathbf{r}|s)$ .  
 134 The difference between these perspectives is illustrated in Figure 3.

### 135 2.3.2 Differing notions of “likelihood”

136 Another major difference is evidenced by divergent usage of the term “likelihood” (Figure 3). In Bayesian  
 137 Encoding, the term “likelihood” is reserved for the abstract relationship between internal model variables  
 138 and sensory data. For instance, one could speak of the “likelihood that this configuration of variables  
 139 in the brain’s model generated the observed image,” or  $p_b(\mathbf{I}|\mathbf{x})$ . This usage supports the idea that the  
 140 quantity being computed is a posterior *over internal variables*. In Bayesian Decoding, on the other hand,  
 141 the “likelihood” refers to a relationship between stimuli and neural responses,  $p(\mathbf{r}|s)$ . This usage supports  
 142 the idea that the quantity of interest is the posterior *over external stimuli*.

### 143 2.3.3 Differing Empirical and Theoretical Motivations

144 Finally, distinguishing Bayesian Encoding and Bayesian Decoding allows one to be more precise on what  
145 data and what normative arguments motivate different theories. Bayesian Decoding can be motivated by  
146 the fact that humans and other species are empirically sensitive to uncertainty and prior experience, as  
147 in the classic psychophysics results on multi-modal cue combination (Ernst and Banks, 2002; Knill and  
148 Pouget, 2004; Alais and Burr, 2004; Körding, 2007; Pouget et al., 2013). The vast literature on optimal or  
149 near-optimal Bayesian perception in controlled tasks motivates the question of how neural circuits facilitate  
150 Bayesian computations *with respect to stimuli in a task*. Bayesian Decoding is further motivated by neural  
151 data which show a correspondence between neural noise, behavioral indications of uncertainty, and decoding  
152 weights in a psychophysics task (Fetsch et al., 2013; Hou et al., 2019; Walker et al., 2019). Importantly, none  
153 of these results constitute direct evidence for inference with respect to an internal model, as hypothesized in  
154 Bayesian Encoding theories.

155 There are three motivations for Bayesian Encoding which are independent of the above motivations  
156 for Bayesian Decoding. First, there is a constraint on *all* well-calibrated statistical models that the prior  
157 must equal the average posterior (Dayan and Abbott, 2001). There is some empirical evidence that this  
158 constraint is satisfied by neural responses in visual cortex (Berkes et al., 2011; Lange and Haefner, 2020).  
159 Second, one can test for signatures of particular inference algorithms and particular internal models trained  
160 on natural stimuli. This approach has been employed by a series of sampling-based inference models and has  
161 successfully reproduced a wide range of neural response properties in early visual cortex (Orbán et al., 2016;  
162 Aitchison and Lengyel, 2016; Echeveste et al., 2019). Third, Bayesian Encoding is often motivated by purely  
163 normative arguments. Any rational agent that faces uncertainty *ought to* compute posterior distributions  
164 over unobserved variables (Jaynes, 2003). However, we emphasize again that existing evidence for near-  
165 optimality in psychophysical tasks only constitutes weak evidence in favor of inference with respect to a  
166 task-independent internal model of the sort usually studied in the Bayesian Encoding literature.

167 While the Encoding and the Decoding perspectives are complementary, it is important to make this  
168 distinction explicit. Failure to do so can lead to confusion and apparently conflicting results on the nature  
169 of the neural code. To illustrate this point, we next construct a model that *encodes* the posterior over  
170 internal variables by sampling and show analytically that it can be exactly *decoded* in a manner consistent  
171 with PPCs. An earlier version of the following section has appeared previously as NeurIPS conference  
172 proceedings (Shivkumar et al., 2018).

## 173 2.4 Decoding Samples from a Linear Gaussian Model is Equivalent to a PPC

174 An earlier version of this example originally appeared in the 2018 NeurIPS conference proceedings (Shivkumar  
175 et al., 2018). At a high level, our example proceeds as follows: we begin with a linear Gaussian internal  
176 generative model and we assume that neurons in V1 approximately infer a posterior distribution over image  
177 features. Inference consists of stochastic samples encoded by spiking responses over time. Next, we expose  
178 this system to stimuli from a task, such as oriented gratings. We then analytically derive the optimal decoder  
179 of task stimuli (e.g. grating orientation) from neural responses, and find that it is a linear PPC. We discuss  
180 a variety of implications, including the connection between neural variability and uncertainty and the role  
181 of nuisance variables in this system.

### 182 2.4.1 Encoding: Neural Sampling in a Linear Gaussian Model

183 We follow previous work in assuming that neurons in primary visual cortex (V1) implement probabilistic  
184 inference in a linear Gaussian model of the input image (Olshausen and Field, 1996, 1997; Hoyer and  
185 Hyvärinen, 2003; Bornschein et al., 2013; Haefner et al., 2016):

$$\mathbf{I} \sim \mathcal{N}(\mathbf{A}\mathbf{x}, \Sigma_{\mathbf{x}}) \quad (2)$$

186 where  $\Sigma_{\mathbf{x}}$  is the covariance of pixel noise in the brain’s generative model. The observed image,  $\mathbf{I}$ , is assumed  
187 to be drawn from a Normal distribution whose mean is a linear combination of “projective fields” ( $\mathbf{PF}_i$ );  
188 the matrix  $\mathbf{A}$  is a feature dictionary with projective fields as its columns:  $\mathbf{A} = (\mathbf{PF}_1, \dots, \mathbf{PF}_n)$ . Each of  
189 the  $n$  projective fields is weighted by a single latent variable,  $\mathbf{x} = (x_1, \dots, x_n)^\top$ , which will later each be



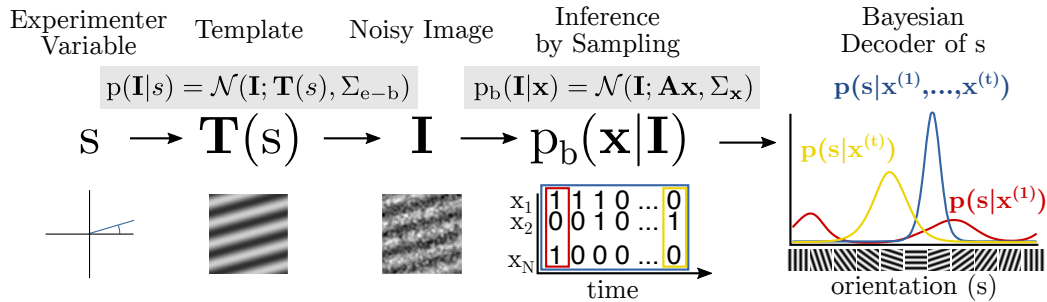


Figure 4: Encoding by sampling followed by decoding of orientation from the samples. Our model performs sampling-based inference over  $\mathbf{x}$  in a probabilistic model of the image,  $\mathbf{I}$ . In a given experiment, the image is generated according to the experimenter’s model that turns a scalar stimulus  $s$ , e.g. orientation, into an image observed by the brain. The samples drawn from the model are then probabilistically “decoded” in order to infer the implied probability distribution over  $s$  from the brain’s perspective. While the samples shown here are binary, our derivation of the PPC is agnostic to whether they are binary or continuous, or to the nature of the brain’s prior over  $\mathbf{x}$ .

190 associated with a single neuron. The main empirical justification for this model consists in the fact that  
 191 under the assumption of a sparse independent prior over the  $\mathbf{x}$ , the model learns projective field parameters  
 192 that resemble the localized, oriented, and bandpass features that characterize V1 neurons when trained on  
 193 natural images (Olshausen and Field, 1996; Bornschein et al., 2013). Hoyer & Hyvarinen (2003) proposed  
 194 that during inference neural responses can be interpreted as samples in such a model. Furthermore, Orban  
 195 et al. (2016) showed that samples from a closely related generative model (Gaussian scale mixture model,  
 196 (Schwartz and Simoncelli, 2001)) could explain many response properties of V1 neurons beyond receptive  
 197 fields. Since our main points are conceptual in nature, we will develop them for the slightly simpler original  
 198 model described above.

199 Given an image,  $\mathbf{I}$ , we assume that neural responses correspond to samples from the posterior distribution,  
 200  $\mathbf{x}^{(t)} \sim p_b(\mathbf{x}|\mathbf{I}) \propto p_b(\mathbf{I}|\mathbf{x})p_b(\mathbf{x})$  where  $p_b(\mathbf{x})$  is the brain’s prior over  $\mathbf{x}$ . The exact form of  $p_b(\mathbf{x})$  will not  
 201 matter for the subsequent decoding arguments. We assume that spikes from a population of  $n$  neurons  
 202 encode instantaneous values of samples from the posterior over  $\mathbf{x}$ , so that each instant, the population  
 203 response,  $\mathbf{x}^{(t)} = (x_1^{(t)}, \dots, x_n^{(t)})^\top$ , represents a sample from the brain’s posterior belief about  $\mathbf{x}|\mathbf{I}$ . Each  
 204 sample of  $x_i$  represents the brain’s instantaneous belief about the intensity of the feature  $\mathbf{PF}_i$  in the image.  
 205 This interpretation is independent of any task demands or assumptions by the experimenter; as discussed  
 206 above,  $\mathbf{x} \rightarrow \mathbf{I}$  is the brain’s *internal* model. In the next section we will show how these samples can also be  
 207 interpreted as a population code over some experimenter-defined quantity like orientation.

#### 208 2.4.2 Decoding: Inferring Task Stimuli from Samples Results in a PPC

209 In many classic neurophysiology experiments, an experimenter presents stimuli that only vary along a scalar  
 210 dimension, such as the orientation of a grating or direction of dot motion (Parker and Newsome, 1998).  
 211 We call this scalar quantity of interest “ $s$ .” We then pose the following decoding question: assuming V1  
 212 implements sampling-based inference as defined in the previous section, what can downstream areas infer  
 213 about  $s$  by observing the sequence of samples produced by V1? An ideal observer would apply Bayes’ rule  
 214 to infer  $p(s|\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}) \propto p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}|s)p(s)$  using knowledge of the *likelihood of generating that set of*  
 215 *samples* for each  $s$ . In the linear Gaussian image model, the optimal decoder can be computed analytically,  
 216 which we do next.

217 We assume the image that is observed by the brain’s sensory periphery (e.g. retinal ganglion cells) is  
 218 defined by a template function  $\mathbf{T}(s)$  plus noise. This template function could, for instance, represent a  
 219 grating of a particular spatial frequency and contrast, or any other shape that is being varied along  $s$  in the  
 220 course of the experiment (Figure 4). We further allow for Gaussian pixel noise around the template  $\mathbf{T}(s)$   
 221 with covariance  $\Sigma_{e-b}$ , which accounts for both (e)xternal pixel noise and noise internal to the (b)rain. This

222 means the likelihood that the brain observes the image  $\mathbf{I}$  conditioned on  $s$  is

$$p(\mathbf{I}|s) = \mathcal{N}(\mathbf{I}; \mathbf{T}(s), \Sigma_{e-b}), \quad (3)$$

223 where  $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Sigma)$  denotes the probability density of a multivariate normal distribution with mean  $\boldsymbol{\mu}$  and  
224 variance  $\Sigma$  evaluated at  $\mathbf{x}$ .

225 With these assumptions, we are able to analytically derive the optimal decoder of  $s$  conditioned on a  
226 sequence of  $t$  independent samples from the posterior,  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\}$ . By Bayes' rule, the optimal decoder  
227 of  $s$  is simply the product of the prior  $p(s)$  with the likelihood of generating those  $t$  samples conditioned on  
228  $s$ . This likelihood term is

$$p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}|s) \propto \mathcal{N}\left(\mathbf{T}(s); \mathbf{A}\bar{\mathbf{x}}_t, \Sigma_{e-b} + \frac{1}{t}\Sigma_{\mathbf{x}}\right) \int \frac{\kappa(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)})}{p_b(\mathbf{I})^t} \mathcal{N}(\mathbf{I}; \boldsymbol{\mu}_{\mathbf{I}}, \Sigma_{\mathbf{I}}) d\mathbf{I}, \quad (4)$$

229 where  $\bar{\mathbf{x}}_t = \frac{1}{t} \sum_{i=1}^t \mathbf{x}^{(i)}$  is the average of all samples up to time  $t$ . A full derivation along, with the exact  
230 form of  $\kappa$ ,  $\boldsymbol{\mu}_{\mathbf{I}}$ , and  $\Sigma_{\mathbf{I}}$  can be found in section S.1 or in Shivkumar et al. (2018). Importantly, as  $t$  gets large,  
231  $\boldsymbol{\mu}_{\mathbf{I}}$  goes to  $\mathbf{A}\bar{\mathbf{x}}_t$ , which means that none of the terms in the integral depend on  $s$ . In the limit of large  $t$ ,  
232 then, the full decoder of  $s$  is given by the much simpler expression,

$$\lim_{t \rightarrow \infty} p(s|\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}) \propto p(s) \mathcal{N}(\mathbf{T}(s); \mathbf{A}\bar{\mathbf{x}}, \Sigma_{e-b}). \quad (5)$$

233 Writing this expression in the canonical form for the exponential family gives

$$\lim_{t \rightarrow \infty} p(s|\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}) \propto g(s) \exp(\mathbf{h}(s)^\top \bar{\mathbf{x}}) \quad \text{where} \quad (6)$$

$$g(s) = \exp\left(-\frac{1}{2}\mathbf{T}(s)^\top \Sigma_{e-b}^{-1} \mathbf{T}(s)\right) p(s) \quad \text{and} \quad (7)$$

$$\mathbf{h}(s) = \mathbf{T}(s)^\top \Sigma_{e-b}^{-1} \mathbf{A}. \quad (8)$$

234 If samples of  $\mathbf{x}$  are encoded by instantaneous neural responses, then firing rates  $\mathbf{r}$  are proportional to  $\bar{\mathbf{x}}$ . We  
235 can then conclude that, in the limit of large  $t$ , this model is equivalent to a linear PPC over  $s$  as defined by  
236 Ma et al. (2006).

### 237 2.4.3 Simulations

238 We simulated this model system estimating the orientation of a grating image, where the generative model  
239 consisted of a mixture of uniformly spaced oriented Gabor patches in the columns of  $\mathbf{A}$ . Figure 5 shows a  
240 numerical simulation of decoded posteriors over  $s$  for different numbers of samples, using the large- $t$  decoder  
241 of equations (6)-(8), to illustrate how drawing additional samples results in a sharper decoded posterior over  
242  $s$ . When only a small number of samples of  $\mathbf{x}$  are drawn, the decoded distributions over  $s$  are both wide  
243 and variable, but get sharper and less variable as the number of samples increases (Figure 5a-c). The black  
244 distribution shown in Figure 5d is both the optimal decoder of  $s$  in the limit of many samples as well as a  
245 PPC over orientation. The bottom row of Figure 5 shows the corresponding spike counts for each neurons  
246 on the  $y$ -axis sorted by the preferred stimulus of each neuron on the  $\mathbf{x}$ -axis.

### 247 2.4.4 The Decoded PPC is Task-Dependent

248 The relationships that we have derived for  $g(s)$  and  $\mathbf{h}(s)$  (equations (7) and (8)) provide insights into the  
249 nature of the PPC that arises in a linear Gaussian model of the inputs. A classic stimulus to consider when  
250 probing and modeling neurons in area V1 is an oriented grating. If the images are identical up to rotation,  
251 and if the prior distribution over orientations is flat, then  $g(s)$  will be constant. Equation (7) shows how  $g(s)$   
252 changes as either of those conditions does not apply, for instance when considering stimuli that vary along  
253 spatial frequency or binocular disparity, rather than orientation, for which the prior significantly deviates  
254 from constant. Further, we can read from equation (8) exactly how the kernels  $\mathbf{h}(s)$ , which characterize  
255 how each neuron contributes to the population code over  $s$ , depend both on the manifold of images defined  
256 by  $\mathbf{T}(s)$ , and on the projective fields contained in the columns of  $\mathbf{A}$ . For an intuition, consider the case of



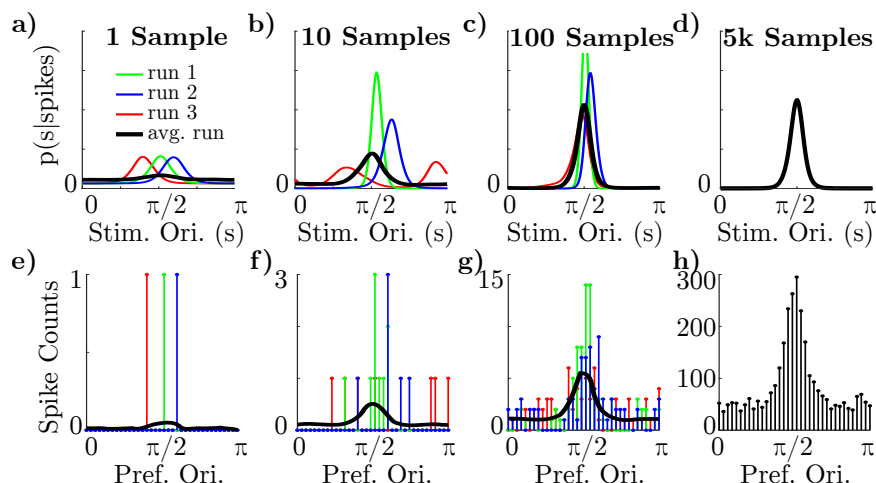


Figure 5: Visualization of the convergence of the decoder after more and more samples. **a-c)** Decoded posterior over  $s$  implied by equation (5) for 1, 10, and 100 samples, respectively. Colored lines are individual sampling runs. Black line is the average posterior over many runs. **d)** Decoded posterior over  $s$  after 5k samples in black, and with the mean  $\bar{x}$  estimated by Variational Bayes in orange. **e-h)** Population responses corresponding to each panel in (a-d) (note different scales from left to right). The three highlighted runs were selected for visualization post-hoc to ensure the first sample contained only a single spike (e), but this is not true in general for all runs.

257 isotropic pixel noise, that is  $\Sigma_{e-b} = \sigma_{e-b}^2 \mathbb{I}$ , in which case  $\mathbf{h}(s)$  is simply the dot product between  $\mathbf{T}(s)$  and  
 258  $\mathbf{PF}_i$  for each neuron, scaled by  $1/\sigma_{e-b}^2$ . The more  $\mathbf{T}(s)^\top \mathbf{PF}_i$  depends on  $s$ , the more informative neuron  $i$ 's  
 259 response is for the posterior over  $s$ .

260 Importantly, the PPC depends as much on the manifold of images defined for a particular experiment,  
 261  $\mathbf{T}(s)$ , as it does on the projective fields of the neurons,  $\mathbf{A}$ . The kernels  $\mathbf{h}(s)$  will be different for gratings  
 262 of different size and spatial frequency, for plaids, or for a house. This is what we mean when we say the  
 263 code over  $s$  is *task-dependent*:  $\mathbf{T}(s)$  is largely arbitrary and up to the experimenter. This means that a  
 264 downstream area forming an estimate of  $s$ , or an area that is combining the information contained in the  
 265 neural responses  $\mathbf{x}$  with that contained in another population (e.g. in the context of cue integration) will  
 266 need to learn the  $\mathbf{h}(s)$  separately for each task.

#### 267 2.4.5 Simultaneous Log- and Direct-Probability Codes

268 One way that questions about the nature of Bayesian inference in the brain has been posed is by considering  
 269 a distinction between Log Probability Codes and Linear or Direct Probability Codes (Barlow, 1969; Pouget  
 270 et al., 2013). Taking the log of equation (6) reveals that the neural responses in our model are linearly related  
 271 to the logarithm of the posterior over  $s$ . By construction, neural responses in our simple model correspond  
 272 to samples, i.e. neither probabilities nor log probabilities over  $\mathbf{x}$ . It is worth noting, however, that samples  
 273 are proportional to probabilities in the special case where all latent variables are binary. In that case, on the  
 274 time scale of a single sample, the response is either 0 or 1, making the firing rate of neuron  $i$  proportional  
 275 to its marginal probability,  $p_b(x_i|\mathbf{I})$ . Such a binary image model has been shown to be as successful as the  
 276 original continuous model of Olshausen & Field (1996) in explaining the properties of V1 receptive fields  
 277 (Henniges et al., 2010; Bornschein et al., 2013), and is supported by studies on plausible implementations of  
 278 sampling in spiking neurons (Buesing et al., 2011; Pecevski et al., 2011). This implies that for the special  
 279 case of binary latents, our neural sampling model is simultaneously a direct probability code (over  $\mathbf{x}_i$ ), and  
 280 a log probability code (over  $s$ ).

#### 281 2.4.6 Dissociating Neural Variability and Uncertainty

282 It is important to appreciate the difference between the brain’s posteriors over  $\mathbf{x}$ , and over  $s$ . The former  
283 represents a belief about an *internal* variable such as the intensity or absence/presence of individual image  
284 elements in the input. The latter represents knowledge about an external stimulus that caused the input  
285 given the neural responses. Neural variability, as modeled here, corresponds to variability in the samples  
286  $\mathbf{x}^{(i)}$  and is directly related to the uncertainty in the posterior over  $\mathbf{x}$ . The uncertainty over  $s$  encoded by the  
287 PPC, on the other hand, depends on the samples only through their *mean, rather than their variance*. Given  
288 sufficiently many samples, the uncertainty over  $s$  is only determined by the noise in the channel between  
289 experimenter and brain ( $\Sigma_{e-b}$ ). This is a sobering point for experiments that seek to determine whether  
290 the brain is sampling by testing the relationship between neural variability and “uncertainty” in broad terms:  
291 in our example model, only uncertainty over  $\mathbf{x}$  but not over  $s$  manifests as neural variability, while  $s$  is the  
292 thing most commonly and naturally manipulated in an experiment.

### 293 3 Discussion

294 Although it is widely agreed that a primary function of sensory neural circuits is to infer *something*, it is  
295 not generally agreed *what* they infer. According to the Bayesian Decoding perspective, neurons represent  
296 distributions over external quantities such as stimuli in a task. According to the Bayesian Encoding per-  
297 spective, neurons represent distributions over variables in an internal model which exists independently of a  
298 task. These are complementary perspectives, and the same system might be interpreted as a fundamentally  
299 different type of distributional code (sampling or a PPC) depending on what variables we assume the system  
300 represents (linear Gaussian features or task stimuli). The question of *how* the brain implements approximate  
301 inference is inextricable from the question of *what* it infers.

302 Historically, sampling-based neural models have taken the Bayesian Encoding approach, asking how neu-  
303 rons could sample from the posterior distribution over variables in an internal model, while PPCs have  
304 primarily been associated with Bayesian Decoding. However, this does not reflect any fundamental distinc-  
305 tion between the two types of distributional codes. Parametric codes can and have been applied to Bayesian  
306 Encoding problems, including both PPCs and other types of parametric codes such as distributed distribu-  
307 tional codes (DDCs) (Vertes and Sahani, 2018). Finally, one could consider cognitive sampling models as a  
308 kind of sampling-based decoding, which have been used to explain a wide variety of perceptual and cognitive  
309 phenomena from multi-stable perception (Gershman et al., 2012) to anchoring and availability biases (Lieder  
310 et al., 2013, 2017). Table 1 provides a list of examples in each of the four categories defined by the sampling  
311 versus parametric and the encoding versus decoding axes.

312 Although Bayesian Decoding is not a trivial problem, it is a weaker form of the Bayesian Brain hypothesis  
313 than Bayesian Encoding. One might call Bayesian Decoding the **Weak** Bayesian Brain Hypothesis, because  
314 it is more descriptive than prescriptive. That is, it describes properties that a neural code ought to have in  
315 order to make the job of downstream circuits “easy,” and it is relatively tractable to ask whether populations  
316 of neurons have those properties – the challenge is to construct  $\mathbf{r}|s$  to realize these properties (Zemel et al.,  
317 1998; Ma et al., 2006). Bayesian Encoding, on the other hand, might be called the **Strong** Bayesian Brain  
318 Hypothesis, because it requires committing to the potentially much harder to falsify idea that the brain  
319 contains an internal generative model of its sensory inputs so that the posterior  $p_b(\mathbf{x}|\mathbf{I})$  is unambiguously  
320 defined.

321 In section 2.3.3, we argued that Bayesian Encoding and Bayesian Decoding have largely disjoint empirical  
322 and theoretical support. Bayesian Decoding can be motivated by the substantial psychophysics literature on  
323 near-optimal perception in the face of ambiguity (Knill and Richards, 1996). However, it would be a mistake  
324 to treat evidence for near-optimal or near-Bayesian behavior in a particular task alone as evidence that  
325 the brain represents probability distributions over variables in an internal model. One could imagine, for  
326 instance, extending our example above to the case where the image features,  $\mathbf{x}$ , are not represented by  
327 samples from their posterior, but by their MAP or mean posterior value. This would be a *point estimate*  
328 over internal variables and thus antithetical to the idea of Bayesian Encoding, but would nonetheless facilitate  
329 many forms of Bayesian Decoding; in fact, neurons encoding only the mean or MAP of  $\mathbf{x}$  in our model would  
330 directly form a linear PPC over  $s$ ! If point estimates of internal model variables are sufficient for Bayesian  
331 Decoding of task quantities, then Bayesian Encoding requires additional justification *outside* the usually-

	Encoding	Decoding
Sampling	Hoyer and Hyvärinen (2003) Berkes et al. (2011) Buesing et al. (2011) Orbán et al. (2016) Haefner et al. (2016) Aitchison and Lengyel (2016) Savin and Denève (2014)*	Lieder et al. (2013) Vul et al. (2014) Gershman et al. (2012) Sanborn and Chater (2016) Lieder et al. (2017)
Parametric	Friston (2005) Beck et al. (2012) Raju and Pitkow (2016) Vertes and Sahani (2018) Zemel et al. (1998) ? Sahani and Dayan (2003) ? Tajima et al. (2016)? Savin and Denève (2014)*	Ma et al. (2006) Beck et al. (2008) Beck et al. (2011) Hou et al. (2019) Zemel et al. (1998) ? Sahani and Dayan (2003) ? Tajima et al. (2016)?

Table 1: Dividing previous work along the lines of sampling versus parametric codes and encoding versus decoding. The fact that there is previous work in all four quadrants emphasizes that these are complementary distinctions. We marked three papers with “?” that are exceptions to the hard division between encoding and decoding. Savin and Denève (2014), marked with “\*”, can similarly be seen as an exception to the hard division between sampling-based and parametric encodings.

332 cited empirical psychophysics literature. The distinction between Bayesian Encoding and Bayesian Decoding  
333 might productively add to the open philosophical question: “if perception is probabilistic, why does it not  
334 seem probabilistic” (Block, 2018; Rahnev et al., 2020).

335 An important question for all Bayesian Encoding models is the extent to which they depend on assump-  
336 tions about the brain’s internal model or inference algorithm. As an example, Berkes et al (2011) compared  
337 the average stimulus-evoked neural activity in visual cortex to spontaneous activity, finding that they be-  
338 come more aligned over the course of development. This is argued to be evidence that the brain develops  
339 an internal statistical model of its sensory inputs in broad terms, since all *well-calibrated* statistical models  
340 have the property that the prior is equal to the average posterior (Dayan and Abbott, 2001). However, this  
341 link requires making crucial assumptions about the nature of the brain’s internal model and its distribu-  
342 tional code. First, Berkes et al assume that neural activity encoding the prior can be directly measured  
343 by recording spontaneous neural activity, i.e. by recording visual cortex in the dark. This assumption is  
344 motivated by the observation that the posterior in scale-mixture models reverts to the prior when contrast  
345 is zero, but is in general not true of other types of image models. As an alternative approach to assuming a  
346 particular type of internal model, one might instead assert that an internal model *exists* while also conceding  
347 that it is *unknown* to us as experimenters. This is the approach taken by Lange & Haefner (2020), who  
348 derived predictions for sensory neural activity from the same principle of learning a *well-calibrated* model,  
349 but without assuming that the brain’s prior can be directly measured.

350 The key step in our example system above which allowed us to interpret samples of  $\mathbf{x}$  as a PPC was  
351 to construct the PPC over a different variable –  $s$ . Still, the distinction between sampling and parametric  
352 codes may also be a false dichotomy *even when considering a single quantity to be inferred*. That is, the  
353 question of whether the brain samples or implements variational inference over its internal  $\mathbf{x}$  may also lead  
354 to a false dichotomy. In principle, it is possible to interpret each sample as implying an entire distribution,  
355 and it is possible to improve variational inference by adding stochasticity to the parameters (Hoffman et al.,  
356 2013). Current proposals for how the brain could implement probabilistic inference are limited by inference  
357 algorithms known from statistics and machine learning, which also tend to divide cleanly into “sampling” or  
358 “variational” methods, but rarely both. One way to advance theories of neural inference, then, may be to  
359 develop statistical algorithms that trade-off the advantages and drawbacks of both sampling and variational  
360 inference (de Freitas et al., 2001; Gershman et al., 2012; Salimans et al., 2015).

## References

- 361  
362 Laurence Aitchison and Máté Lengyel. The Hamiltonian Brain: Efficient Probabilistic Inference with  
363 Excitatory-Inhibitory Neural Circuit Dynamics. *PLoS Computational Biology*, pages 1–24, 2016.
- 364 David Alais and David Burr. The Ventriloquist Effect Results from Near-Optimal Bimodal Integration.  
365 *Current Biology*, 14(3):257–262, 2004.
- 366 Mihály Bányai, Andreea Lazar, Liane Klein, Johanna Klon-Lipok, Marcell Stippinger, Wolf Singer, and  
367 Gergő Orbán. Stimulus complexity shapes response correlations in primary visual cortex. *Proceedings of*  
368 *the National Academy of Sciences*, 116(7):2723–2732, 2019.
- 369 H. B. Barlow. Pattern Recognition and the Responses of Sensory Neurons. *Annals of the New York Academy*  
370 *of Sciences*, 156(2):872–881, 1969.
- 371 Jeffrey M. Beck, Wei Ji Ma, Roozbeh Kiani, Timothy D. Hanks, Anne K. Churchland, Jamie Roitman,  
372 Michael N. Shadlen, Peter E. Latham, and Alexandre Pouget. Probabilistic Population Codes for Bayesian  
373 Decision Making. *Neuron*, 36(6):1142–1152, 2008.
- 374 Jeffrey M. Beck, Peter E. Latham, and Alexandre Pouget. Marginalization in Neural Circuits with Divisive  
375 Normalization. *J. Neurosci.*, 31(43):15310–15319, 2011.
- 376 Jeffrey M. Beck, Katherine Heller, and Alexandre Pouget. Complex Inference in Neural Circuits with  
377 Probabilistic Population Codes and Topic Models. *Advances in Neural Information Processing Systems*,  
378 25:3068–3076, 2012.
- 379 Pietro Berkes, Gergo Orbán, Máté Lengyel, and József Fiser. Spontaneous Cortical Activity Reveals Hall-  
380 marks of an Optimal Internal Model of the Environment. *Science*, 331(January):83–87, 2011.
- 381 Christopher M Bishop. *Pattern Recognition and Machine Learning*. Springer, Cambridge, 2006.
- 382 Ned Block. If perception is probabilistic, why does it not seem probabilistic? *Philosophical Transactions of*  
383 *the Royal Society B: Biological Sciences*, 373(1755), 2018.
- 384 Jörg Bornschein, Marc Henniges, and Jörg Lücke. Are V1 Simple Cells Optimized for Visual Occlusions? A  
385 Comparative Study. *PLoS Computational Biology*, 9(6), 2013.
- 386 Lars Buesing, Johannes Bill, Bernhard Nessler, and Wolfgang Maass. Neural dynamics as sampling: A model  
387 for stochastic computation in recurrent networks of spiking neurons. *PLoS Computational Biology*, 7(11),  
388 2011.
- 389 Peter Dayan and Larry F. Abbott. *Theoretical Neuroscience: Computational and Mathematical Modeling of*  
390 *Neural Systems*. MIT Press, London, 2001.
- 391 Nando de Freitas, Pedro Højen-Sørensen, Michael I. Jordan, and Stuart Russel. Variational MCMC. *UAI*,  
392 2001.
- 393 Rodrigo Echeveste, Laurence Aitchison, Guillaume Hennequin, and Máté Lengyel. Cortical-like dynamics in  
394 recurrent circuits optimized for sampling-based probabilistic inference. *bioRxiv*, page 696088, 2019.
- 395 Marc O Ernst and Martin S Banks. Humans integrate visual and haptic information in a statistically optimal  
396 fashion. *Nature*, 415(6870):429–433, 2002.
- 397 Christopher R. Fetsch, Alexandre Pouget, Gregory C. DeAngelis, and Dora E. Angelaki. Neural correlates of  
398 reliability-based cue weighting during multisensory integration. *Nature Neuroscience*, 15(1):146–54, 2011.
- 399 Christopher R. Fetsch, Gregory C. DeAngelis, and Dora E. Angelaki. Bridging the gap between theories of  
400 sensory cue integration and the physiology of multisensory neurons. *Nature Reviews Neuroscience*, 14(6):  
401 429–442, 2013.

- 402 József Fiser, Pietro Berkes, Gergo Orbán, and Máté Lengyel. Statistically optimal perception and learning:  
403 from behavior to neural representations. *Trends in cognitive sciences*, 14(3):119–30, mar 2010.
- 404 Karl J. Friston. A theory of cortical responses. *Philosophical transactions of the Royal Society of London.*  
405 *Series B*, 360:815–836, 2005.
- 406 Samuel J. Gershman and Jeffrey M. Beck. Complex Probabilistic Inference: From Cognition to Neural  
407 Computation. In Ahmed Moustafa, editor, *Computational Models of Brain and Behavior*, chapter Complex  
408 Pr. Wiley-Blackwell, 2016.
- 409 Samuel J. Gershman, Edward Vul, and Joshua B. Tenenbaum. Multistability and perceptual inference.  
410 *Neural Computation*, 24(1):1–24, 2012.
- 411 Ralf M. Haefner, Pietro Berkes, and József Fiser. Perceptual Decision-Making as Probabilistic Inference by  
412 Neural Sampling. *Neuron*, 90:649–660, 2016.
- 413 Guillaume Hennequin, Yashar Ahmadian, Daniel B Rubin, Máté Lengyel, and Kenneth D Miller. The Dy-  
414 namical Regime of Sensory Cortex: Stable Dynamics around a Single Stimulus-Tuned Attractor Account  
415 for Patterns of Noise Variability. *Neuron*, 98:846–860, 2018.
- 416 Marc Henniges, Gervasio Puertas, Jörg Bornschein, Julian Eggert, and Jörg Lücke. Binary Sparse Coding.  
417 In Vincent Vigneron, Vicente Zarzoso, Eric Moreau, Rémi Gribonval, and Emmanuel Vincent, editors,  
418 *Latent Variable Analysis and Signal Separation*, pages 450–457, 2010.
- 419 Matthew D. Hoffman, David M. Blei, Chong Wang, and John Paisley. Stochastic variational inference.  
420 *Journal of Machine Learning Research*, 14:1303–1347, 2013.
- 421 Han Hou, Qihao Zheng, Yuchen Zhao, Alexandre Pouget, and Yong Gu. Neural Correlates of Optimal  
422 Multisensory Decision Making under Time-Varying Reliabilities with an Invariant Linear Probabilistic  
423 Population Code. *Neuron*, 104:1–12, 2019.
- 424 Patrik O. Hoyer and Aapo Hyvärinen. Interpreting neural response variability as monte carlo sampling of  
425 the posterior. *Advances in Neural Information Processing Systems*, 17(1):293–300, 2003.
- 426 E. T. Jaynes. *Probability Theory: The Logic of Science*. Cambridge University Press, New York, 2003.
- 427 David C. Knill and Alexandre Pouget. The Bayesian brain: the role of uncertainty in neural coding and  
428 computation. *Trends in Neurosciences*, 27(12):712–9, dec 2004.
- 429 David C. Knill and Whitman Richards, editors. *Perception as Bayesian Inference*. Cambridge University  
430 Press, New York, NY, 1996.
- 431 Konrad P Körding. Decision Theory: What "Should" the Nervous System Do? *Science Review*, 318, 2007.
- 432 Richard D. Lange and Ralf M. Haefner. Task-induced neural covariability as a signature of approximate  
433 Bayesian learning and inference. *bioRxiv*, 2020.
- 434 Falk Lieder, Thomas L. Griffiths, and Noah D. Goodman. Burn-in , bias , and the rationality of anchoring.  
435 *Advances in Neural Information Processing Systems*, 25:1–9, 2013.
- 436 Falk Lieder, Thomas L. Griffiths, Quentin J M Huys, and Noah D. Goodman. Empirical Evidence for  
437 Resource-Rational Anchoring and Adjustment. *Psychonomic Bulletin & Review*, 2017.
- 438 Wei Ji Ma, Jeffrey M. Beck, Peter E. Latham, and Alexandre Pouget. Bayesian inference with probabilistic  
439 population codes. *Nature Neuroscience*, 9(11):1432–1438, 2006.
- 440 Kevin P. Murphy. *Machine Learning: A Probabilistic Perspective*. The MIT Press, Cambridge, MA, 2012.
- 441 Bruno A. Olshausen and David J. Field. Emergence of simple-cell receptive field properties by learning a  
442 sparse code for natural images, 1996.



- 443 Bruno a Olshausen and David J. Field. Sparse coding with an incomplete basis set: a strategy employed by  
444 V1?, 1997.
- 445 Gergő Orbán, Pietro Berkes, József Fiser, and Máté Lengyel. Neural Variability and Sampling-Based Prob-  
446 abilistic Representations in the Visual Cortex. *Neuron*, 92(2):530–543, 2016.
- 447 A J Parker and William T. Newsome. Sense and the Single Neuron: Probing the Physiology of Perception.  
448 *Annual Review of Neuroscience of neuroscience*, 21:227–277, 1998.
- 449 Dejan Pecevski, Lars Buesing, and Wolfgang Maass. Probabilistic inferences general graphical models  
450 through sampling in stochastic networks of spiking neurons. *PLOS Computational Biology*, 7(12), 2011.
- 451 Alexandre Pouget, Jeffrey M. Beck, Wei Ji Ma, and Peter E. Latham. Probabilistic brains: knowns and  
452 unknowns. *Nature Neuroscience*, 16(9):1170–8, 2013.
- 453 Dobromir Rahnev, Ned Block, Janneke Jehee, and Rachel Denison. Is perception probabilistic? In *Cognitive*  
454 *Computational Neuroscience*, 2020.
- 455 Rajkumar V. Raju and Xaq Pitkow. Inference by Reparameterization in Neural Population Codes. *Advances*  
456 *in Neural Information Processing Systems*, 30, 2016.
- 457 Maneesh Sahani and Peter Dayan. Doubly Distributional Population Codes: Simultaneous Representation  
458 of Uncertainty and Multiplicity. *Neural Computation*, 15:2255–2279, 2003.
- 459 Tim Salimans, Diederik P. Kingma, and Max Welling. Markov Chain Monte Carlo and Variational Inference:  
460 Bridging the Gap. In *International Conference on Machine Learning*, pages 1218–1226, 2015.
- 461 Adam N Sanborn. Types of approximation for probabilistic cognition: Sampling and variational. *Brain and*  
462 *Cognition*, 2015.
- 463 Adam N Sanborn and Nick Chater. Bayesian Brains without Probabilities. *Trends in Cognitive Sciences*, 20  
464 (12):883–893, 2016.
- 465 Cristina Savin and Sophie Denève. Spatio-temporal representations of uncertainty in spiking neural networks.  
466 *Advances in Neural Information Processing Systems*, 2014.
- 467 Odelia Schwartz and Eero P Simoncelli. Natural signal statistics and sensory gain control. *Nature Neuro-*  
468 *science*, 4(8):819–825, 2001.
- 469 Sabyasachi Shivkumar, Richard D. Lange, Ankani Chattoraj, and Ralf M. Haefner. A probabilistic population  
470 code based on neural samples. *Advances in Neural Information Processing Systems*, 31:7070—7079, 2018.
- 471 Chihiro I. Tajima, Satohiro Tajima, Kowa Koida, Hidehiko Komatsu, Kazuyuki Aihara, and Hideyuki Suzuki.  
472 Population code dynamics in categorical perception. *Nature Scientific Reports*, 6(22536):1–13, 2016.
- 473 Eszter Vertes and Maneesh Sahani. Flexible and accurate inference and learning for deep generative models.  
474 *Advances in Neural Information Processing Systems*, 31, 2018.
- 475 Hermann von Helmholtz. *Treatise on Physiological Optics*. The Optical Society of America, 1925.
- 476 Edward Vul, Noah D. Goodman, Thomas L. Griffiths, and Joshua B. Tenenbaum. One and done? Optimal  
477 decisions from very few samples. *Cognitive Science*, 38(4):599–637, 2014.
- 478 Martin J. Wainwright and Michael I. Jordan. Graphical Models, Exponential Families, and Variational  
479 Inference. *Foundations and Trends in Machine Learning*, 1(1–2):1–305, 2008.
- 480 Edgar Y Walker, R. James Cotton, Wei Ji Ma, and Andreas S Tolias. A neural basis of probabilistic  
481 computation in visual cortex. *Nature Neuroscience*, 23:122–129, 2019.
- 482 Richard S. Zemel, Peter Dayan, and Alexandre Pouget. Probabilistic Interpretation of Population Codes.  
483 *Neural Computation*, 10(2):403–430, 1998.



## 484 S Supplemental Text

### 485 S.1 Derivation of the optimal decoder from samples

486 Here we derive a slightly more general result than is stated in the main text by considering arbitrary  
 487 covariance matrices: we consider here the case where  $\mathbf{I}$  is distributed with mean  $\mathbf{T}(s)$  and covariance  $\Sigma_{e-b}$ ,  
 488 and the brain's internal model generates images with mean  $\mathbf{Ax}$  and covariance  $\Sigma_{\mathbf{x}}$ . The probability of  
 489 drawing a single neural sample,  $\mathbf{x}^{(i)}$ , given an observed image  $\mathbf{I}$ , by assumption, equal to the posterior  
 490 probability of  $\mathbf{x}$  in the brain's internal model. The probability of drawing a sequence of  $t$  independent  
 491 samples of  $\mathbf{x}$  is,<sup>1</sup>

$$\begin{aligned} p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)} | \mathbf{I}) &= \prod_{i=1}^t p(\mathbf{x}^{(i)} | \mathbf{I}) \\ &= \prod_{i=1}^t \frac{p_b(\mathbf{I} | \mathbf{x}^{(i)}) p_b(\mathbf{x}^{(i)})}{p_b(\mathbf{I})} \\ &= \frac{1}{p_b(\mathbf{I})^t} \prod_{i=1}^t p_b(\mathbf{I} | \mathbf{x}^{(i)}) p_b(\mathbf{x}^{(i)}). \end{aligned}$$

Our results primarily follow from this identity for the product of two multivariate normal distributions:

$$\begin{aligned} \mathcal{N}(\mathbf{y}; \mu_1, \Sigma_1) \mathcal{N}(\mathbf{y}; \mu_2, \Sigma_2) &= \mathcal{N}(\mathbf{y}; \mu_3, \Sigma_3) \mathcal{N}(\mu_1; \mu_2, \Sigma_1 + \Sigma_2) \quad (\text{S1}) \\ \Sigma_3 &= (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1} \\ \mu_3 &= \Sigma_3 (\Sigma_1^{-1} \mu_1 + \Sigma_2^{-1} \mu_2) \end{aligned}$$

492 Letting  $\bar{\mathbf{x}}_{t'} = \frac{1}{t'} \sum_{i=1}^{t'} \mathbf{x}^{(i)}$  denote the running mean of the samples up to  $t'$ , it follows from the above product  
 493 identity that

$$\prod_{i=1}^t \underbrace{\mathcal{N}(\mathbf{I}; \mathbf{Ax}^{(i)}, \Sigma_{\mathbf{x}})}_{p_b(\mathbf{I} | \mathbf{x}^{(i)})} = \mathcal{N}\left(\mathbf{I}; \mathbf{A}\bar{\mathbf{x}}_t, \frac{1}{t} \Sigma_{\mathbf{x}}\right) \prod_{t'=2}^t \mathcal{N}(\mathbf{Ax}^{(t')}; \mathbf{A}\bar{\mathbf{x}}_{t'-1}, \frac{t'}{t'-1} \Sigma_{\mathbf{x}}). \quad (\text{S2})$$

We next absorb all terms that do not depend on  $s$  or  $\mathbf{I}$  into  $\kappa(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)})$ . Specifically, let

$$\kappa(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}) = \prod_{i=1}^t p_b(\mathbf{x}^{(i)}) \prod_{t'=2}^t \mathcal{N}(\mathbf{Ax}^{(t')}; \mathbf{A}\bar{\mathbf{x}}_{t'-1}, \frac{t'}{t'-1} \Sigma_{\mathbf{x}}).$$

After simplifying further, this can be written in terms of a ratio of Gaussian densities with mean zero, times the product of priors on each  $\mathbf{x}$ :

$$\kappa(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}) = \frac{\prod_{i=1}^t \mathcal{N}(\mathbf{Ax}^{(i)}; \mathbf{0}, \Sigma_{\mathbf{x}}) p_b(\mathbf{x}^{(i)})}{\mathcal{N}(\mathbf{A}\bar{\mathbf{x}}_t; \mathbf{0}, \frac{1}{t} \Sigma_{\mathbf{x}})}.$$

494 Then, the likelihood of drawing a particular set of  $t$  independent samples of  $\mathbf{x}$  conditioned on  $\mathbf{I}$  is

$$p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)} | \mathbf{I}) \propto \frac{\kappa(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)})}{p_b(\mathbf{I})^t} \mathcal{N}\left(\mathbf{I}; \mathbf{A}\bar{\mathbf{x}}_t, \frac{1}{t} \Sigma_{\mathbf{x}}\right). \quad (\text{S3})$$

495 .

496 Since a decoder looking only at samples of  $\mathbf{x}$  has no direct access to the image, the likelihood for a  
 497 full sequence of samples conditioned on  $s$  requires marginalizing over all possible images  $\mathbf{I}$  that could be  
 498 generated conditioned on a fixed  $s$ :

$$p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)} | s) = \int p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)} | \mathbf{I}) p(\mathbf{I} | s) d\mathbf{I}.$$

<sup>1</sup>We write  $p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)} | \mathbf{I})$  rather than  $p_b(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)} | \mathbf{I})$  because while  $\mathbf{x}$  is part of the brain's internal model, the *samples* of  $\mathbf{x}$  are not, but are viewed through the lens of an outside observer or optimal decoder.

499 Substituting in (S3), this is

$$p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)} | s) = \int \frac{\kappa(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)})}{p_b(\mathbf{I})^t} \mathcal{N}\left(\mathbf{I}; \mathbf{A}\bar{\mathbf{x}}, \frac{1}{t}\Sigma_{\mathbf{x}}\right) p(\mathbf{I} | s) d\mathbf{I}.$$

500 Next, making use of the assumption that  $\mathbf{I} | s$  is a multivariate normal centered on  $\mathbf{T}(s)$  with pixel covariance  
501  $\Sigma_{e-b}$  and applying the multivariate normal product identity (S1), it follows that

$$\begin{aligned} p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)} | s) &= \int \frac{\kappa(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)})}{p_b(\mathbf{I})^t} \mathcal{N}\left(\mathbf{I}; \mathbf{A}\bar{\mathbf{x}}_t, \frac{1}{t}\Sigma_{\mathbf{x}}\right) \mathcal{N}(\mathbf{I}; \mathbf{T}(s), \Sigma_{e-b}) d\mathbf{I} \\ &= \mathcal{N}\left(\mathbf{T}(s); \mathbf{A}\bar{\mathbf{x}}_t, \Sigma_{e-b} + \frac{1}{t}\Sigma_{\mathbf{x}}\right) \int \frac{\kappa(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)})}{p_b(\mathbf{I})^t} \mathcal{N}(\mathbf{I}; \boldsymbol{\mu}_{\mathbf{I}}, \Sigma_{\mathbf{I}}) d\mathbf{I}, \end{aligned} \quad (4 \text{ restated})$$

502 where

$$\begin{aligned} \Sigma_{\mathbf{I}} &= (t\Sigma_{\mathbf{x}}^{-1} + \Sigma_{e-b}^{-1})^{-1} \\ \boldsymbol{\mu}_{\mathbf{I}} &= \Sigma_{\mathbf{I}} (\Sigma_{e-b}^{-1} \mathbf{T}(s) + t\Sigma_{\mathbf{x}}^{-1} \mathbf{A}\bar{\mathbf{x}}_t). \end{aligned}$$

503 As we will show below, the first term in (4),  $\mathcal{N}(\mathbf{T}(s); \mathbf{A}\bar{\mathbf{x}}_t, \dots)$ , implies that the decoder is a linear PPC.  
504 The integral in (4) requires further discussion. First, note that as the number of samples,  $t$ , increases,  $\Sigma_{\mathbf{I}}$   
505 shrinks towards zero, and  $\boldsymbol{\mu}_{\mathbf{I}}$  goes to  $\mathbf{A}\bar{\mathbf{x}}_t$ , which implies that  $\mathcal{N}(\mathbf{I}; \boldsymbol{\mu}_{\mathbf{I}}, \Sigma_{\mathbf{I}})$  goes to a delta distribution around  
506  $\mathbf{A}\bar{\mathbf{x}}$ . This implies that for large  $t$ , the integral ceases to depend on  $s$ , and hence can be ignored by a decoder.  
507 Thus, for large  $t$ , we have

$$p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)} | s) \propto \mathcal{N}(\mathbf{T}(s); \mathbf{A}\bar{\mathbf{x}}_t, \Sigma_{e-b}), \quad (S4)$$

508 where the proportionality should be understood in the context of decoding  $s$ , and is only approximate for  
509 finite  $t$ . Note that when  $t$  is small, it may still be the case that the integral in (4) does not depend strongly  
510 on  $s$ . This is the case, for instance, if the brain's internal model assigns equal probability to all  $\mathbf{T}(s)$ , in  
511 which case  $p_b(\mathbf{I})$  evaluated at  $\boldsymbol{\mu}_{\mathbf{I}}$  does not depend on  $s$ .

512 Applying Bayes' rule to *decode*  $s$  from the samples of  $\mathbf{x}$ , and absorbing all terms that do not contain  $s$   
513 into the proportionality constant, (S4) implies

$$p(s | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}) \propto p(s) \mathcal{N}(\mathbf{T}(s); \mathbf{A}\bar{\mathbf{x}}, \Sigma_{e-b}). \quad (S5)$$

We can now rewrite this expression in the canonical form for the exponential family

$$p(s | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}) \propto g(s) \exp(\mathbf{h}(s)^\top \bar{\mathbf{x}}) \quad \text{where} \quad (6 \text{ restated})$$

$$g(s) = \exp\left(-\frac{1}{2} \mathbf{T}(s)^\top \Sigma_{e-b}^{-1} \mathbf{T}(s)\right) p(s) \quad \text{and} \quad (7 \text{ restated})$$

$$\mathbf{h}(s) = \mathbf{T}(s)^\top \Sigma_{e-b}^{-1} \mathbf{A}. \quad (8 \text{ restated})$$

514 Equating samples of  $\mathbf{x}$  with instantaneous neural responses, the firing rate  $\mathbf{r}$  is proportional to  $\bar{\mathbf{x}}$ . We can  
515 then conclude that, in the limit of large  $t$ , this model is a linear PPC over  $s$  as defined by (Ma et al., 2006).