

Bayesian Encoding and Decoding as Distinct Perspectives on Neural Coding

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

Abstract

One of the most influential, and controversial, ideas in neuroscience has been to understand the brain in terms of Bayesian computations. Unstated differences in how Bayesian ideas are operationalized across different models make it difficult to ascertain both which empirical data support which models, and how Bayesian computations might be implemented by neural circuits. In this paper, we make one such difference explicit by identifying two distinct philosophies that underlie existing neural models of Bayesian inference: one in which the brain recovers experimenter-defined structures in the world from sensory neural activity (Decoding), and another in which the brain represents latent quantities in an internal model that explains its inputs (Encoding). These philosophies require profoundly different assumptions about the nature of inference in the brain, and lead to different interpretations of empirical data. Here, we characterize and contrast both philosophies in terms of motivations, empirical support, and relationship to neural data. We also show that this implicit difference in philosophy underlies some of the debate on whether neural activity is better described as a sampling-based, or a parametric, distributional code. Using a simple model of primary visual cortex as an example, we show mathematically that it is possible that the very same neural activity can be described as probabilistic inference by neural sampling in the *Encoding* framework while also forming a linear probabilistic population code (PPC) in the *Decoding* framework. This demonstrates that certain families of Encoding and Decoding models are compatible with each other rather than competing explanations of data. In sum, Bayesian Encoding and Bayesian Decoding are distinct, non-exclusive philosophies, and appreciating their similarities and differences will help organize future work and allow for stronger empirical tests about the nature of inference in the brain.

1 Introduction

According to the Bayesian Brain hypothesis, neural circuits carry out statistical computations by combining prior knowledge with new evidence, combining multiple sources of information according to their reliability, and taking actions that account for uncertainty. 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 evolution. 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 and perceptual decision-making (Knill and Richards, 1996; Kersten et al., 2004; Fiser et al., 2010; 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 computational standpoint, sampling and parametric codes each 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 used to explain spontaneous

*equal contribution

20 cortical activity (Berkes et al., 2011), neural variability (Hoyer and Hyvärinen, 2003; Orbán et al., 2016; Festa
21 et al., 2021), structure in noise correlations (Haefner et al., 2016; Bányai et al., 2019), and onset transients
22 and oscillations (Aitchison and Lengyel, 2016; Hennequin et al., 2018; Echeveste et al., 2020). Meanwhile,
23 parametric codes have been cited in explanations of contrast-invariant tuning curves (Ma et al., 2006), near-
24 linearity during cue-combination (Fetsch et al., 2011, 2013), evidence integration dynamics in parietal cortex
25 (Beck et al., 2008; Hou et al., 2019), divisive normalization (Beck et al., 2011), and more (Pouget et al., 2013).
26 Importantly, sampling and parametric codes have so far always been discussed and compared as competing
27 and mutually exclusive mathematical models of the same neural circuits, with no decisive evidence presented
28 favoring one over the other model. Notably, contrast-invariant tuning and divisive normalization have also
29 been replicated by sampling models (Orbán et al., 2016; Echeveste et al., 2020).

30 The primary goal of this paper is to characterize and contrast two distinct perspective on the Bayesian
31 Brain hypothesis, which we call **Bayesian Encoding** and **Bayesian Decoding**. These are complementary
32 perspectives that make different assumptions about the nature of the inference problems faced by the brain,
33 and are supported or falsified by different kinds of empirical data. We argue that not making their differences
34 explicit has led to confusion about how to interpret empirical data. In particular, we describe how the above
35 debate on whether neural responses are better modeled as samples or parameters is complicated by the
36 fact that sampling codes usually make assumptions consistent with Bayesian Encoding while parametric
37 codes often make assumptions consistent with Bayesian Decoding. However, neither the connection between
38 Bayesian Encoding and sampling, nor between Bayesian Decoding and parametric codes, is a necessary
39 consequence of either theory. Indeed, there are Encoding models built on parametric codes, Decoding
40 models based on sampling, and still other models that contain elements of both approaches (Section 2.4 and
41 Table 1 below).

42 Finally, we illustrate the complementary nature of these two philosophies using a simple model of primary
43 visual cortex (Shivkumar et al., 2018). In this example, we construct a sampling-based *Encoding* model
44 based on a linear Gaussian model of natural images (Olshausen and Field, 1996, 1997), and derive the
45 implied *Decoding* model. We show that firing rates in this model form a canonical kind of parametric code:
46 a Probabilistic Population Code (PPC). There is thus no inherent contradiction in saying that the brain is
47 *both* sampling (in the “Bayesian Encoding” sense) *and* represents parameters (in the “Bayesian Decoding”
48 sense), and depending on the encoding models’ generative model, and the considered task, this parametric
49 code may even be a linear PPC. We conclude with a discussion of distributional neural codes in general.

50 2 Bayesian Encoding vs Bayesian Decoding

51 We follow the seminal work of Zemel et al. (1998) in assuming that patterns of neural activity represent entire
52 probability distributions over a variable, not just a point estimate of it, i.e. that they form a *distributional*
53 code. The nature of this “variable” and its relationship to neural response is key to the distinction between
54 the Bayesian Encoding and the Bayesian Decoding frameworks.

55 2.1 Bayesian Encoding

56 We define **Bayesian Encoding** as the view that there exists a probability distribution over some quantity of
57 potential interest to the brain, and that the primary function of sensory neurons is to compute and represent
58 an approximation to this distribution. We use the term “encoding” because the probability distribution
59 that neurons are hypothesized to represent conceptually precedes the actual neural responses. That is, in
60 Bayesian encoding models, there exists a *reference distribution* that is defined independently of how neurons
61 actually respond, and which is approximately encoded by neural responses.

62 Bayesian Encoding requires a source for the reference distribution. In the context of the sensory system,
63 this typically takes the form of an internal generative model of sensory inputs, and the distribution to be
64 encoded is the posterior over latent variables in that model (Figure 1a-b). With this perspective, the goal
65 of sensory areas of the brain is to learn a statistical model of its sensory inputs (Dayan et al., 1995; Dayan
66 and Abbott, 2001; Fiser et al., 2010; Berkes et al., 2011) in which sensory observations, such as an image
67 on the retina, are explained as the result of higher order causes. Whereas the information on the retina
68 is highly mixed – objects, lights, textures, and optics interact in complex ways to create an image – the
69 internal model aims to explain sensory data in terms of unobserved causes that are often assumed to be

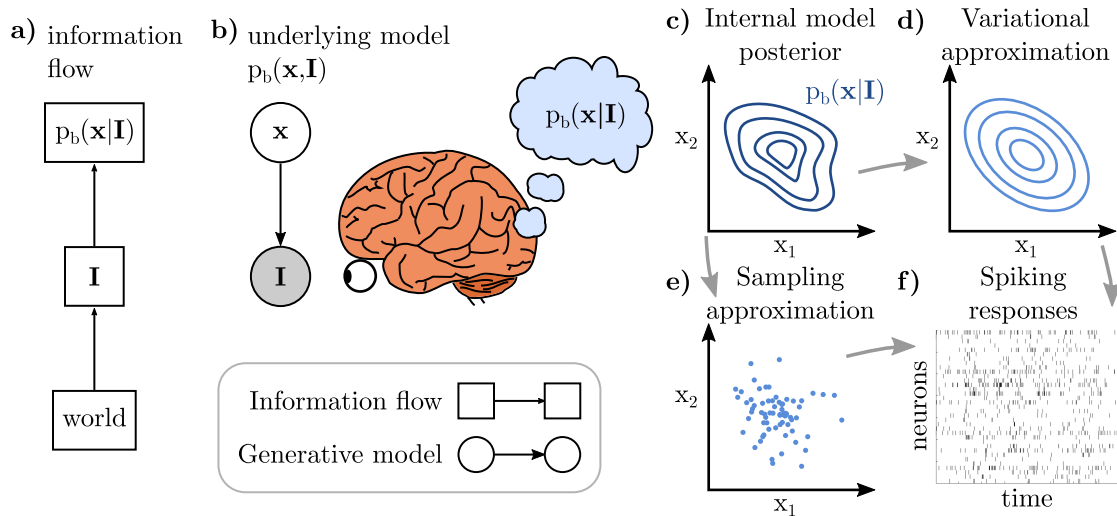


Figure 1: Visualization of Bayesian Encoding. **a)** Diagram of information flow: the world provides sensory inputs (\mathbf{I}), which then give rise to inferences about latent variables (\mathbf{x}). **b)** Bayesian Encoding typically assumes that the brain has an internal model of its inputs, and that perceptual inferences are about variables in this internal model, not necessarily corresponding to quantities in the external world *per se*. With Bayesian Encoding, it is also typical to assume that the internal model is *task-independent* and that the brain always computes a posterior over internal variables, $p_b(\mathbf{x}|\mathbf{I})$, regardless of whether \mathbf{I} is a highly controlled stimulus in a task or encountered in the wild. **c-f)** The defining feature of Bayesian Encoding is the existence of a reference distribution (c), typically the posterior over a set of latent variables, \mathbf{x} , given a sensory measurement, \mathbf{I} . One then assumes an approximation scheme such as variational inference (c→d) or sampling (c→e), and that this approximation is then realized in patterns of neural activity (f).

70 sparse and independent (von Helmholtz, 1925; Olshausen and Field, 1996; Bell and Sejnowski, 1997). A
71 generative model makes this process explicit by assigning prior probabilities to the (co)occurrence of causes
72 (represented by latent variables) and by quantifying the likelihood of a particular configuration of the causes
73 for generating a particular sensory observation. The encoded posterior distribution in this framework is
74 defined over the latent variables in this statistical model.

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

$$76 \quad p_b(\mathbf{x}|\mathbf{I}) = \frac{p_b(\mathbf{I}|\mathbf{x})p_b(\mathbf{x})}{p_b(\mathbf{I})}. \quad (1)$$

77 We use the subscript b in $p_b(\cdot)$ to refer to quantities in the brain’s internal model, and to distinguish them
78 from other types of probabilities such as a decoder’s uncertainty. A prototypical case of Bayesian Encoding
79 poses the question of how neural circuits could compute and represent the posterior distribution $p_b(\mathbf{x}|\mathbf{I})$
80 for any sensory \mathbf{I} , given the internal model that the brain has learned (Figure 1c), and how it can learn
81 this internal model in the first place. In general, exact inference is an intractable problem (Murphy, 2012;
82 Wainwright and Jordan, 2008; Bishop, 2006), leading to the question of how the brain could compute and
83 represent an *approximation* to the true posterior (Figure 1d-f), and what the nature of this approximation is.
84 This line of reasoning motivates work on “neurally plausible approximate inference algorithms,” including
85 approaches with connections to sampling-based inference (Figure 1e), as well as approaches inspired by
86 variational inference techniques, related to parametric neural codes (Figure 1d) (reviewed in Fiser et al.
87 (2010); Sanborn (2015); Gershman and Beck (2016)).

88 2.2 Bayesian Decoding

89 We define **Bayesian Decoding** as the perspective in which neural activity is treated as *given*, and emphasis
90 is placed on the statistical uncertainty of a decoder observing those neural responses. Bayesian Decoding is
91 closely related to ideal observer models in psychophysics, involving tasks that require the estimation of scalar
92 aspects of a presented stimulus (e.g. its orientation or its contrast) or a decision whether the stimulus belongs
93 to one of two or more discrete classes (e.g. “left” or “right”). Of course, any stimulus s that elicits neural
94 responses \mathbf{r} is optimally decoded by computing $p(s|\mathbf{r})$ (Figure 2). The key question within the Bayesian
95 Decoding framework is this: what conditions must the stimulus-driven neural activity ($p(\mathbf{r}|s)$) fulfill such
96 that the decoder ($p(s|\mathbf{r})$) is *simple*, e.g. linear and invariant to nuisance? For instance, imposing linearity
97 and invariance constraints on the decoder implies constraints on tuning curves and the distribution of neural
98 noise (Zemel et al., 1998; Ma et al., 2006).

99 Bayesian Decoding is closely related to familiar notions of optimal neural decoding. Classically, decoding
100 is either a tool for assessing information content in neural responses or a mechanistic model of how they
101 impact behavior. In the Bayesian setting, the emphasis is on how neural activity is interpreted by the rest
102 of the brain and influences behavior, and how this depends on the brain’s uncertainty about a behaviorally-
103 relevant stimulus.

104 Probabilistic Population Codes (PPCs), as introduced by Ma et al (2006), exemplify the Bayesian De-
105 coding approach. PPCs construct a Bayesian decoder that is both simple and invariant to nuisance: if
106 a population of neurons tuned to s has “Poisson-like” variability, then the optimal decoder is part of the
107 exponential family with firing rates as natural parameters. This is a particularly convenient representation
108 for taking products of two distributions as required by cue-integration (Ma et al., 2006) and evidence ac-
109 cumulation Beck et al. (2008). Equally important is the notion of *invariance* afforded by a PPC: as long
110 as nuisance variables such as image contrast or dot coherence only multiplicatively scale tuning curves, the
111 decoder can ignore them.

112 Importantly, under the assumption that the brain employs a computationally convenient neural code,
113 linearity for cue combination and multiplicative gain by nuisance variables become *predictions* of PPCs.
114 In classical decoding approaches, neural responses are simply “given,” not prescribed by a theory. In the
115 Bayesian Decoding framework generally, and in the case of PPCs in particular, imposing constraints on the
116 decoder constrains the possible set of evoked response distributions, $p(\mathbf{r}|s)$. These constraints have then
117 been formulated as predictions and tested empirically (Fetsch et al., 2011, 2013; Pouget et al., 2013; Hou
118 et al., 2019).

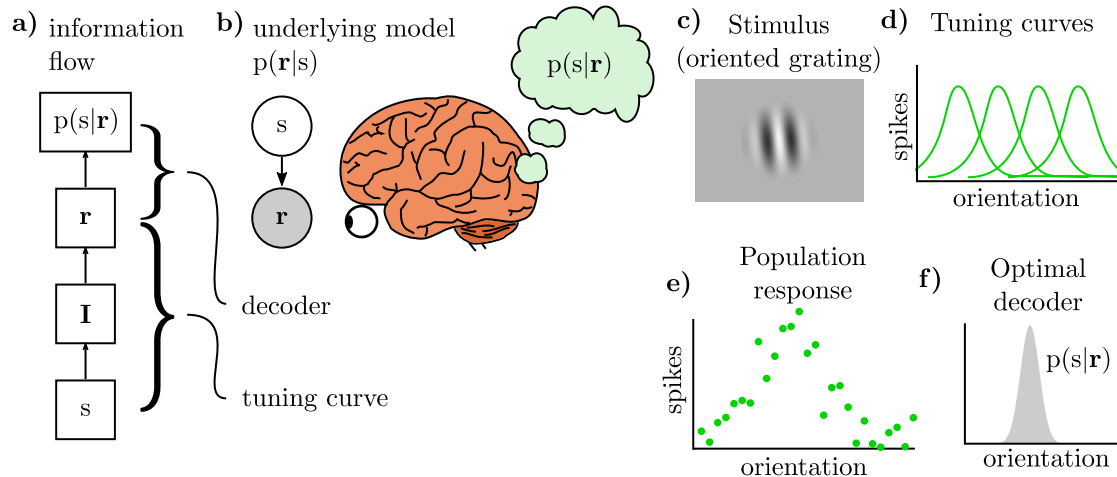


Figure 2: Visualization of Bayesian Decoding. **a)** Diagram of information flow: a quantity of interest (s) in the world elicits neural responses (r), mediated by sensory inputs (I). The decoding question is how the brain forms an internal estimate, \hat{s} , from r . **b)** The underlying probabilistic model assumes that s generates r , so inference is the problem of recovering s from r . **c)** The decoding problem usually begins with a stimulus, such as the orientation of a grating of a given spatial frequency, size, location, and contrast. **d-f)** Given a population of neurons’ tuning curves to s (d) and an observation of spikes on a single trial (e), an optimal decoder computes $p(s|r)$ (f).

119 2.3 Contrasting Bayesian Encoding and Bayesian Decoding

120 There are four key differences between the Bayesian Encoding and Bayesian Decoding perspectives, which we
 121 will discuss in each of the following sections: (1) what they assume the brain is inferring, (2) what the terms
 122 “likelihood” and “posterior” refer to, (3) the role of neural responses in the theory, and (4) the empirical
 123 data and other arguments used to motivate them. As our goal is to summarize and categorize a large and
 124 diverse sub-field, there will be exceptions to each rule, but we expect these distinctions to be useful for
 125 framing further discussions.

126 2.3.1 Differences in what is assumed to be inferred

127 An integral part of the Bayesian Encoding framework is the existence of an abstract internal model that
 128 could in principle be implemented *in silico* or in the brains of other individuals or other species. Deriving
 129 predictions for neural data requires an additional linking hypothesis on the nature of distributional codes,
 130 such as whether neurons sample or encode variational parameters, and how either samples or parameters
 131 correspond to observable biophysical quantities like membrane potentials and spike times or spike counts.
 132 Bayesian Encoding thus decomposes the question of what sensory neurons compute into two parts: first,
 133 what is the internal model which defines optimal inference (the reference distribution), and second, how do
 134 neural circuits carry out approximate inference in that model (e.g. sampling or parametric)?

135 The brain’s internal model is typically assumed to have been calibrated through exposure to natural
 136 stimuli (Dayan et al., 1995; Dayan and Abbott, 2001; Berkes et al., 2011) and to only change slowly with
 137 exposure to new stimuli in adult brains. For this reason, the generative model in Bayesian Encoding models,
 138 especially in the case of early sensory areas, is often assumed to be independent of experimental context.
 139 For instance, if the brain’s internal model comprises patches of local image features, then it is assumed that
 140 the brain infers and encodes the same set of image features, whether viewing natural scenes or artificial
 141 stimuli in a task (Haefner et al., 2016; Orbán et al., 2016; Shivkumar et al., 2018; Bányai et al., 2019). The
 142 assumption of calibration in a Bayesian Encoding framework also makes predictions for how the internal
 143 model should change in response to the statistics of sensory inputs during development (Berkes et al., 2011),
 144 and to extensive exposure to stimuli in a particular task (Lange and Haefner, 2022).

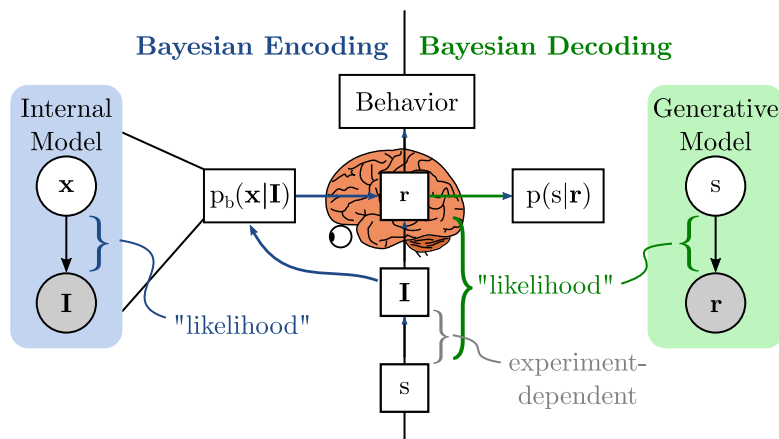


Figure 3: Side-by-side comparison of Bayesian Encoding and Bayesian Decoding. In both frameworks, it is understood that there exists a mechanistic, biophysical connection between stimuli (\mathbf{I}), sensory neural responses (\mathbf{r}), and behavior. In the Bayesian Decoding framework, emphasis is placed on the uncertainty of a decoder estimating a 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} can be seen as an instance of *downward causation* between levels of abstraction, where changes to the posterior (at the algorithmic level) imply changes to neural responses (at the implementation level) (Campbell, 1974; Yablo, 1992; Lange and Haefner, 2022). 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} . In any psychophysical task, the link between s and \mathbf{I} depends on the experiment (gray bracket). Importantly, this means that the “likelihood” in a Bayesian Decoding model depends on choices made by the experimenter (such as their choice of stimuli), but not in a Bayesian Encoding model.

145 In contrast, Bayesian Decoding models are typically applied in the context of estimating task-relevant
146 variables. For instance, in a motion discrimination task, a Bayesian Decoding question would be how the
147 brain represents uncertainty over directions of motion. Importantly, answering this question does not require
148 a generative model of all possible naturally-occurring motion stimuli, nor does it require a true or correct
149 reference distribution over stimuli; it requires only a statistical model of the relation between scalar motion
150 direction in a particular task (and possibly nuisance variables like coherence) and neural responses, $p(\mathbf{r}|s)$,
151 making it observable experimentally. The difference between (typical realizations of) the Bayesian Encoding
152 and Decoding perspectives is illustrated in Figure 3.

153 We emphasize that Bayesian Encoding *typically* but not *necessarily* involves a task-independent internal
154 generative model, and Bayesian Decoding likewise is *typically* but not *necessarily* applied to task-specific
155 variables. In fact, where the ideas of encoding and decoding distributions first appeared in Zemel et al (1998),
156 the encoded distribution was over task quantities (such as \mathbf{x} being motion or heading direction), without
157 specifying an internal generative model, and the decoding problem was framed as the inverse to the encoding
158 problem – that is, recovering the encoded $p(\mathbf{x}|\dots)$ from \mathbf{r} . This again emphasizes the complementary
159 nature of these philosophies: we are free to apply the Bayesian Decoding framework to variables in a task-
160 independent internal model (given \mathbf{r} , what do we know about \mathbf{x} or $p_b(\mathbf{x}|\mathbf{I})?$), or to apply the logic of Bayesian
161 Encoding to task-specific quantities (construct \mathbf{r} to encode a desired $p(s|\mathbf{I})$), but such examples are rare. In
162 the remainder of this paper, we will associate Bayesian Encoding with task-independent internal generative
163 models and Bayesian Decoding with variables in a task, and in the Discussion we will return to the possibility
164 that task variables are explicitly represented as part of the brain’s internal model.

165 2.3.2 Differing notions of likelihood

166 Another difference in philosophy is evidenced by divergent usage of the term “likelihood” (Figure 3). In the
167 typical Bayesian Encoding setting, the term “likelihood” is reserved for the abstract relationship between
168 internal model variables and sensory observations. For instance, one could speak of the “likelihood that
169 this configuration of variables in the brain’s model generated the observed image,” or $p_b(\mathbf{I}|\mathbf{x})$. This usage
170 supports the idea that the quantity being computed is a posterior *over variables in a generative model of*
171 *sensory data*. In the typical Bayesian Decoding setting, on the other hand, the “likelihood” refers to a
172 relationship between stimuli and neural responses, $p(\mathbf{r}|s)$. This usage supports the idea that the quantity of
173 interest is the posterior *over external stimuli in a task*.

174 2.3.3 Differences in the relationship between distributions and neural activity

175 Bayesian Encoding models require two distinct assumptions: first, what is the source of the reference dis-
176 tribution to be encoded (e.g. what is the brain’s internal model $p_b(\mathbf{x}, \mathbf{I})$); and second, what is the linking
177 hypothesis that maps probability distributions to neural activity (Figure 1)? This approach of starting with
178 the encoded distribution abstracts away from the details of neural circuits that must actually *implement*
179 inference. Take the model of Orbán et al. (2016) for example. In this work, the authors assume that neurons
180 in primary visual cortex implement a sampling algorithm to encode the posterior distribution over latent
181 variables in a Gaussian Scale Mixture model. This specifies the reference distribution. It is then assumed
182 that, by some *unspecified* mechanism, the trajectory of a set of neurons’ membrane potentials over time
183 traces out real-valued samples from the posterior, and that these membrane potentials elicit spikes through
184 a nonlinear accumulation process. This specifies the linking hypothesis, or the map from the reference dis-
185 tribution to neural data. This model successfully reproduced a diverse set of known properties about V1
186 (Orbán et al., 2016), but it is not a mechanistic model. From a modeling standpoint, the way that an input
187 image elicits neural activity is *mediated* by the reference posterior: an example of “downward causation”
188 (Campbell, 1974; Yablo, 1992).

189 While for Encoding models there is a clear separation of computational model and neural link, they still of
190 course beg the question of how inference in the computational model is implemented in neural circuits. Prior
191 work has investigated the question of how biologically-plausible recurrent circuits could implement sampling
192 (Bill et al., 2015; Probst et al., 2015; Aitchison and Lengyel, 2016; Petrovici et al., 2016; Dold et al., 2019;
193 Echeveste et al., 2020) or message-passing (George and Hawkins, 2009; Beck et al., 2012; Raju and Pitkow,
194 2016; Grabska-Barwinska et al., 2013; Grabska-Barwińska et al., 2017; George et al., 2018) through their

195 dynamics. However, in these examples there is typically a cost to increased biological plausibility, either by
196 degrading the quality of the encoded distribution, or by degrading the match to empirical neural data.

197 Bayesian Decoding models, in contrast, do not distinguish between the uncertainty in an underlying
198 probabilistic model and the uncertainty of a downstream brain area applying a Bayesian decoder to some
199 neural activity. As a result, Decoding models replace the assumption about the link to neural activity with
200 a *constraint* on the relationship between stimuli (s) and neural activity (\mathbf{r}).

201 To illustrate this point, let us revisit one of the motivating examples for distributional codes of Zemel
202 et al. and contrast the Encoding and Decoding approaches. Consider a rat who is placed into a water maze
203 and must navigate to a hidden platform (Morris, 1984). Initially, the rat may be uncertain about which
204 direction it is facing, e.g. if opposite walls of the maze look the same its correct belief about direction will
205 be bimodal. Similar to the orientation of a grating, head direction is a scalar variable in $[0, 2\pi]$ that we will
206 call s . In the Encoding approach, one might begin by asking what is $p_b(s|\mathbf{I})$ according to an internal model
207 of the environment, where \mathbf{I} stands for the sensory cues the rat uses to orient itself. The distribution $p_b(s|\mathbf{I})$
208 determines how uncertain the rat *ought* to be, according to the internal model. Continuing the Encoding
209 approach, one would then adopt a linking hypothesis (sampling, parametric, etc.) whereby $p_b(s|\mathbf{I})$ is encoded
210 in neural activity \mathbf{r} . In an Encoding model, the encoding of a distribution may be imperfect and lossy, or it
211 may contain more information about the distribution than is being used by downstream circuits. In either
212 case, the way a downstream circuit *uses* the neural activity will generally differ from a Bayesian decoder.

213 Applying the Bayesian Decoding framework to the same problem, we would say that the uncertainty
214 in $p(s|\mathbf{r})$ is the primary kind of uncertainty we should be concerned with, and that there is no distinction
215 between this and the rat’s internal model. Crucially, this does not trivialize representations of uncertainty as
216 “just” a matter of optimal decoding. In the Decoding approach, there may still be an ideal uncertainty that
217 the rat *ought* to have when it is first placed into the maze; however the assumption is that this uncertainty
218 is realized through the way \mathbf{r} is tuned to its inputs \mathbf{I} . That is, it is left to the brain (evolution, learning)
219 to have carefully constructed tuning functions $p(\mathbf{r}|\mathbf{I})$, such that $p(s|\mathbf{r})$ is equal to $p_b(s|\mathbf{I})$ (Ma et al., 2006).
220 One way that the Encoding and Decoding perspectives can become identical, then, is when the decoded
221 distribution $p(s|\mathbf{r})$ equals the reference distribution $p_b(s|\mathbf{I})$. From the Encoding point of view, this requires
222 that the encoding of $p_b(s|\mathbf{I})$ into \mathbf{r} is lossless (or “efficient” in the terminology of Beck et al. (2012)). From
223 the Decoding point of view, they are identical by assumption.

224 Finally, the preceding discussion points to an important practical difference between Encoding and De-
225 coding philosophies in terms of how neural responses are interpreted by downstream areas. In a Decoding
226 model, a downstream area implicitly applies Bayes’ rule to the neural responses arriving from an upstream
227 area to extract information about a stimulus. In an Encoding model, on the other hand, upstream neural
228 activity represents samples or parameters that are then processed by the downstream area according to
229 an underlying approximate inference algorithm, which generally will *not* apply Bayes’ rule to the incoming
230 activity directly. To put it another way, if one assumes that upstream neural activity encodes samples or
231 parameters in an approximate inference algorithm, then there is an important conceptual difference between
232 a downstream area that interprets upstream activity *as samples* or *as parameters* (as in Encoding models),
233 and a downstream area that *decodes* the activity it receives by applying Bayes’ rule to the neural activity.

234 2.3.4 Differing Empirical and Theoretical Motivations

235 Finally, distinguishing Bayesian Encoding and Bayesian Decoding allows one to be more precise on what
236 data and what normative arguments motivate different theories. Bayesian Decoding can be motivated by
237 the fact that humans and other species are empirically sensitive to uncertainty and prior experience, as
238 in the classic psychophysics results on multi-modal cue combination (Ernst and Banks, 2002; Knill and
239 Pouget, 2004; Alais and Burr, 2004; Körding, 2007; Angelaki et al., 2009; Pouget et al., 2013). The large
240 literature on optimal or near-optimal Bayesian perception in controlled tasks motivates the question of how
241 neural circuits facilitate Bayesian computations *with respect to stimuli in a task*, which are often scalar or
242 low-dimensional. With the additional assumption that the neural representation of task-relevant aspects
243 of stimuli is formatted to be easily decoded (e.g. linear and invariant to nuisance (Ma et al., 2006)), this
244 line of reasoning has given rise to predictions for neural data. These predictions have since been largely
245 confirmed for the representation of self-motion in dorsal medial superior temporal area (MSTd) (Fetsch
246 et al., 2011, 2013; Hou et al., 2019). Bayesian Decoding is further motivated by experimental data showing

	Bayesian Encoding	Bayesian Decoding
Sampling-based representation	Hoyer and Hyvärinen (2003) Pecevski et al. (2011) Berkes et al. (2011) Buesing et al. (2011) Gershman et al. (2012) Savin and Denève (2014) Probst et al. (2015) Orbán et al. (2016) Haefner et al. (2016) Aitchison and Lengyel (2016) Festa et al. (2021)	Moreno-Bote et al. (2011) [†]
Parametric representation	Zemel et al. (1998) Sahani and Dayan (2003) Friston (2005) George and Hawkins (2009) Beck et al. (2012) Raju and Pitkow (2016) Vertes and Sahani (2018) Tajima et al. (2016)*	Ma et al. (2006) Beck et al. (2008) Beck et al. (2011) Hou et al. (2019) Tajima et al. (2016)* Moreno-Bote et al. (2011) [†]

Table 1: Classifying previous work on Bayesian neural models according to whether they construct Bayesian Encoding or Decoding models, and whether they use a sampling-based or a parametric neural representation. Tajima et al., marked with “*” contains elements of both encoding and decoding. Moreno-Bote et al., marked with “†”, contains elements of both sampling-based and parametric decoding.

247 a correspondence between non-parametric likelihood functions, neural noise, and behavioral indications of
 248 uncertainty (Walker et al., 2019).

249 Importantly, none of these results constitute direct evidence for inference with respect to an (usually
 250 high-dimensional) internal model of natural stimuli, as hypothesized in typical Bayesian Encoding theories
 251 (Rahnev, 2019; Koblinger et al., 2021). The three lines of support for Bayesian Encoding models are largely
 252 independent of the above motivations for Bayesian Decoding. First, Bayesian Encoding can be motivated by
 253 the purely normative argument that any rational agent that faces uncertainty *ought to* compute probability
 254 distributions over unobserved variables, as long as those variables directly enter into calculations of expected
 255 utility (Jaynes, 2003). Second, there is some empirical evidence for a key prediction of Bayesian encoding
 256 models: a general constraint on *all* well-calibrated statistical models is that the prior must equal the av-
 257 erage posterior (Dayan and Abbott, 2001). Existing observations suggest that this constraint is satisfied
 258 in early visual cortex, as evidenced by changes in neural responses in primary visual cortex over develop-
 259 ment (Berkes et al., 2011) and task-learning (Haefner et al., 2016; Lange and Haefner, 2022). Third, there
 260 is empirical evidence for signatures of particular inference algorithms and particular internal models fit to
 261 natural stimuli. This approach has been employed by a series of sampling-based inference models and has
 262 successfully reproduced a wide range of neural response properties in early visual cortex (Orbán et al., 2016;
 263 Aitchison and Lengyel, 2016; Echeveste et al., 2020). A similar approach has also been taken by parametric
 264 models, where neural circuits have been hypothesized to implement the dynamics of a variational inference
 265 algorithm (Friston, 2005; George and Hawkins, 2009; Beck et al., 2012; Grabska-Barwinska et al., 2013; Raju
 266 and Pitkow, 2016; George et al., 2018; Lavin et al., 2018). We emphasize again that existing evidence for
 267 Bayesian-like behavior in psychophysical tasks only constitutes weak evidence in support of the idea that
 268 the brain computes distributions over variables in a task-independent internal model, as usually studied in
 269 the Bayesian Encoding literature (Rahnev, 2019; Koblinger et al., 2021).

270 2.4 Classification of existing models

271 Historically, sampling-based neural models have taken the Bayesian Encoding approach, asking how neurons
 272 could sample from the posterior distribution over variables in an internal model, while PPCs have primarily

273 been studied in the context of inference of low-dimensional task-relevant quantities. However, this does not
274 reflect a fundamental distinction between the two types of distributional codes. Parametric codes can and
275 have been used in Bayesian Encoding models to approximate the posterior over variables in a generative
276 model, including Probabilistic Population Codes (PPCs) (Beck et al., 2012; Grabska-Barwinska et al., 2013;
277 Raju and Pitkow, 2016), Distributed Distributional Codes (DDCs) (Vertes and Sahani, 2018), and others
278 (Friston, 2005; George and Hawkins, 2009; Lavin et al., 2018; George et al., 2018). Markov Chain Monte
279 Carlo (MCMC) sampling has been used to explain perceptual bistability (Moreno-Bote et al., 2011; Gershman
280 et al., 2012), which could be seen as a form of sampling-based Bayesian Decoding (cf. Hohwy et al. (2008)).
281 To summarize, Table 1 provides a list of examples in each of the four categories defined by the sampling
282 versus parametric and the encoding versus decoding axes. The fact that there is previous work in all four
283 quadrants emphasizes that these are complementary distinctions.

284 2.5 Case Study: primary visual cortex (V1)

285 We now focus on primary visual cortex (V1) to provide a concrete example illustrating and further elaborating
286 on our general points above. Focusing on area V1 has the advantage that much neurophysiological data exists,
287 and both encoding and decoding approaches have enjoyed some success. We will first briefly describe existing
288 work from both perspectives, and then use a simple example to show how they can lead to very different
289 conclusions about the neural code. To that end, we will assume a Bayesian Encoding model that encodes the
290 posterior over internal variables by sampling and show analytically how to derive the corresponding Bayesian
291 Decoding model, obtaining a parametric representation (PPC) (Shivkumar et al., 2018).

292 2.5.1 Bayesian Encoding models for V1

293 The starting point for the Bayesian Encoding approach, applied to V1, is an assumption about the brain’s
294 generative model $p_b(\mathbf{x}, \mathbf{I})$. That is, we must specify what is \mathbf{x} , the variable assumed to be inferred and
295 represented by V1 neurons, and how \mathbf{x} is related to the sensory observations, \mathbf{I} . For simple cells in area V1,
296 Olshausen and Field proposed a linear Gaussian likelihood $\mathbf{I} \sim \mathcal{N}(\mathbf{A}\mathbf{x}, \Sigma_{\mathbf{I}})$ with a sparse independent prior
297 $p_b(\mathbf{x})$ as the brain’s internal model Olshausen and Field (1996, 1997). (We use the notation $\mathbf{I} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$
298 to indicate a random variable drawn from a multivariate normal distribution, and $\mathcal{N}(\mathbf{I}; \boldsymbol{\mu}, \Sigma)$ to denote its
299 density function). In this model, the observed retinal image, \mathbf{I} , is assumed to be a linear combination of
300 “projective fields” (\mathbf{PF}_i) plus unexplained pixel noise $\Sigma_{\mathbf{I}}$; the matrix \mathbf{A} is a feature dictionary with projective
301 fields as its columns: $\mathbf{A} = (\mathbf{PF}_1, \dots, \mathbf{PF}_n)$. Each of the n projective fields is weighted by a single latent
302 variable, $\mathbf{x} = (x_1, \dots, x_n)^\top$. Intuitively, in this model, V1 activity is assumed to somehow represent beliefs
303 about what values for \mathbf{x} best explain a given retinal image, \mathbf{I} .

304 The next assumption in the Encoding framework is about the neural code, or how the posterior distribu-
305 tion, $p_b(\mathbf{x}|\mathbf{I})$, is represented by neural responses, \mathbf{r} . Continuing the previous example, Olshausen and Field
306 assumed that each x_i was represented by a single neuron whose firing rate was proportional to the most
307 probable value for x_i given an image (maximum a posteriori, MAP): $r_i \propto \operatorname{argmax}_{x_i} p_b(\mathbf{x}|\mathbf{I})$. In this model, a
308 single neuron represents the most likely intensity with which a visual feature is present in the image. This is
309 not a fully Bayesian Encoding model in the sense that only the MAP, but not the full posterior distribution
310 $p_b(\mathbf{x}|\mathbf{I})$ is encoded in neural responses. Empirical support for this model is based on the observation that
311 learning (fitting) this model on natural images yields visual features (\mathbf{PF}_i) that are localized, oriented, and
312 band-pass filtered, implying neural responses and receptive fields with similar properties – just as observed
313 empirically (Olshausen and Field, 1996).

314 Subsequent work has both modified and extended this generative model, and combined it with different
315 neural codes. Hoyer and Hyvärinen proposed that neural responses can be understood as samples from the
316 posterior in the same generative model to qualitatively explain the variability and mean-variance relationship
317 of neural responses. Schwartz and Simoncelli extended the generative model to a Gaussian scale mixture
318 model to explain the empirically observed contrast normalization of V1 responses, and Orbán et al. (2016)
319 found agreement between the predictions of a Gaussian scale mixture model combined with neural sampling
320 and a wide range of observations related to the stimulus-dependence of neural variability. Bornschein et al.
321 (2013) proposed a variation of the generative model of Olshausen and Field using a nonlinear Gaussian
322 likelihood and/or binary as opposed to continuous latents \mathbf{x} , and Coen-Cagli et al. (2015) found that a

323 further extension to the generative model in the form of a Mixture of Gaussian Scale Mixture model could
324 explain center-surround interactions in V1. Finally, Haefner et al. (2016) combined the generative model of
325 Olshausen and Field with the ideal observer model of a discrimination task to explain choice probabilities
326 and task-dependent noise correlations of V1 neurons.

327 The key shared element of all these models is an explicit assumption about the computational variable
328 \mathbf{x} that is being represented, and how this variable is related to the sensory observations \mathbf{I} . This model being
329 adapted to natural inputs is an important constraint, and the model parameters are usually obtained by
330 fitting the model to sets of natural images. These models are then general purpose and can be queried
331 using natural inputs or images presented in a task. The fundamental framing is how V1 neurons *encode* the
332 posterior over \mathbf{x} given an arbitrary input, \mathbf{I} .

333 2.5.2 Bayesian Decoding models for V1

334 The starting point for the Bayesian Decoding approach, applied to V1, is a measurement of the conditional
335 probability (or likelihood of s), $p(\mathbf{r}|s)$, for some stimulus s to which V1 neurons are tuned, and that is
336 hypothesized to be represented, such as orientation. Importantly, this means that the measured likelihood
337 is to some extent under experimental control, since the experimenter chooses what images correspond to
338 each value of s (e.g. the size, contrast, or spatial frequency of a grating). In general, for an arbitrary
339 s , this likelihood will be very complicated reflecting the fact that s cannot easily be decoded from \mathbf{r} (e.g.
340 object identity from V1). However, for V1 responses it has empirically been found that the optimal Bayesian
341 decoder for orientation is approximately linear in spike counts and invariant to contrast, a classic nuisance
342 variable (Graf et al., 2011). This finding has been interpreted as meaning that V1 activity “represents”
343 orientation. In conjunction with the Poisson-like neural response variability in V1, this implies that the
344 beliefs of a Bayesian decoder of orientation applied to the neural responses are part of the exponential
345 family. Furthermore, the sufficient statistics are linear in the neural responses. Such a neural representation
346 of a belief has been called a Probabilistic Population Code (PPCs) as introduced by Ma et al..

347 The same logic applies to other candidates for s that modulate the responses of V1 neurons in a straight-
348 forward manner, such as spatial frequency or location. The key element of the Bayesian Decoding approach
349 is taking the perspective of downstream circuits trying to extract information about s from V1 activity: how
350 is the information about s formatted in V1 activity, and is $p(s|\mathbf{r})$ “simple”? In contrast to the Bayesian
351 Encoding perspective which justifies its choice of \mathbf{x} by its fit to natural images and its ability to predict
352 neural responses, the Bayesian Decoding perspective justifies its choice of s by desirable properties of an
353 efficient decoder, e.g. linearity and invariance to nuisance variables.

354 2.5.3 Example model where decoding a stimulus s from encoded samples results in a PPC

355 Since our main points are conceptual in nature, we will develop the link between the Encoding and the
356 Decoding approach for the simple case of a linear Gaussian model with a Gaussian prior, under the assump-
357 tion of a sampling-based neural code. These simplifying assumptions make the difference between Encoding
358 and Decoding clear and analytically tractable, but are not meant to maximize biological plausibility. For
359 instance, the posterior variance in this model is independent of \mathbf{I} , whereas it is well-known that neural
360 response variance is stimulus-dependent, and this effect is captured by neural sampling models with less
361 trivial generative models (Orbán et al., 2016; Bányai et al., 2019; Festa et al., 2019). Beginning with a more
362 complicated Encoding model would lead to a more complicated relationship to Decoding models (e.g. where
363 the Decoder is more complex, e.g. a nonlinear PPC, or not even in the exponential family). Importantly,
364 the core of our argument remains: that an Encoding model based on one type of neural code (e.g. sampling)
365 and a Decoding model based on another type (e.g. parametric) need not be in contradiction with each other,
366 and offer complementary perspectives on the same system.

367 Given an image, \mathbf{I} , we assume that V1 neurons *encode* the posterior $p_b(\mathbf{x}|\mathbf{I})$ by sampling t values from
368 from the posterior distribution, $\mathbf{x}^{(t)} \sim p_b(\mathbf{x}|\mathbf{I}) \propto p_b(\mathbf{I}|\mathbf{x})p_b(\mathbf{x})$ where $p_b(\mathbf{x})$ is the brain’s prior over \mathbf{x} (Hoyer
369 and Hyvärinen, 2003). We assume that responses from a population of n neurons correspond to samples
370 from the posterior over \mathbf{x} , so that at each instant, the population response, $\mathbf{r}^{(t)}$, equals the sample $\mathbf{x}^{(t)}$. Each
371 sample of x_i (or r_i) represents the brain’s instantaneous belief about the intensity of the feature \mathbf{PF}_i in the
372 image.

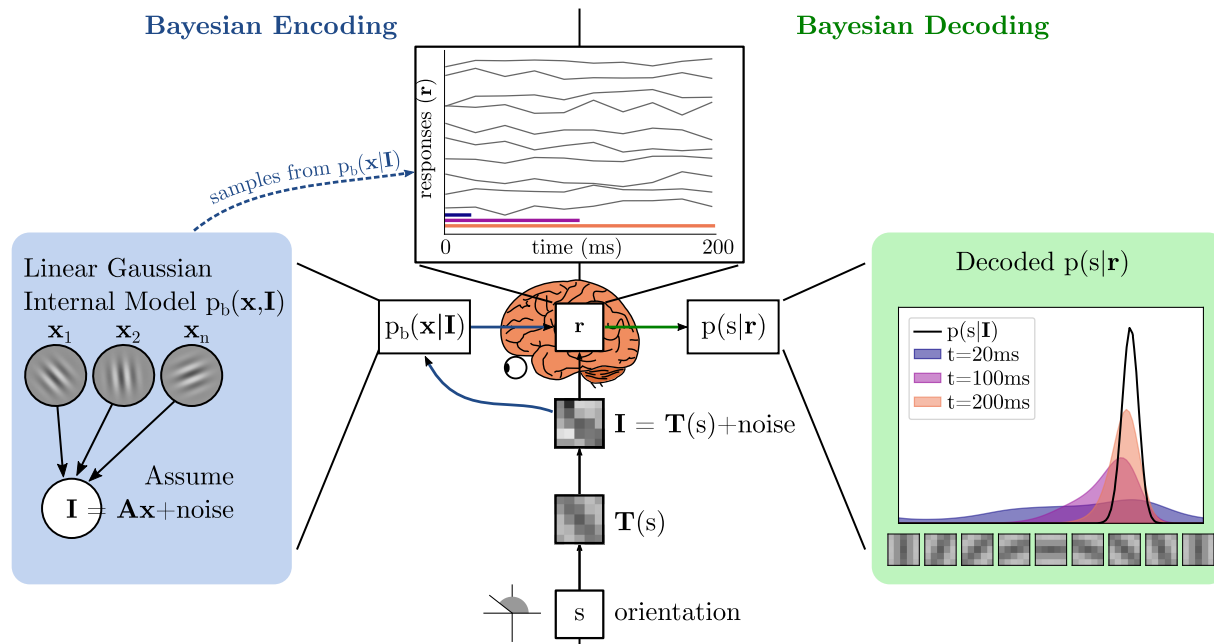


Figure 4: Encoding by sampling followed by decoding of orientation from the samples in a simplified model. As in Figure 3, Encoding elements are on the left, and Decoding elements are on the right. In our example model, the brain performs sampling-based inference over \mathbf{x} in a probabilistic model of images, here a Linear Gaussian model. In a given experiment, the image is generated according to an experimenter-defined process that turns a scalar stimulus s , e.g. orientation, into an image observed by the brain. To simplify, neural responses \mathbf{r} are assumed to reflect instantaneous real-valued samples of \mathbf{x} drawn from the posterior $p_b(\mathbf{x}|\mathbf{I})$. In our simulation we drew 10 samples and assumed 20ms per sample. The samples drawn from the model are then probabilistically “decoded” to a probability distribution over s . This distribution sharpens as more samples are observed. The optimal decoder for any t is a linear PPC.

We will now apply the Bayesian Decoding approach to the sequence of samples produced by the sampling-based Encoding model described above. An ideal observer applies Bayes’ rule to infer $p(s|\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(t)})$ using knowledge of the probabilistic relationship between samples (\mathbf{x} or \mathbf{r}) and s :

$$\begin{aligned}
 p(s|\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(t)}) &\propto p(s) p_b(\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(t)}|s) \\
 &\propto p(s) \int p(\mathbf{I}|s) p_b(\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(t)}|\mathbf{I}) d\mathbf{I}.
 \end{aligned}
 \tag{2}$$

373 That is, the optimal decoder combines knowledge of (i) how likely an image \mathbf{I} is to generate a set of samples
 374 of \mathbf{x} (or \mathbf{r}), and (ii) how likely a stimulus value s is to generate an image \mathbf{I} . In general, this decoded
 375 distribution over s may be arbitrarily complex and intractable. One factor that is under experimental
 376 control is the “template” function $\mathbf{T}(s)$ which renders an image, such as a grating with orientation s . This
 377 provides the link between s and \mathbf{I} in equation (2). In our model, we assume that the input the brain receives
 378 is a noisy version of that template (Figure 4).

The first simplification to the general form of the optimal decoder in (2) we can derive, under the assumption of a Gaussian likelihood, is to show that the posterior over s depends only on the mean rate of \mathbf{r} (i.e. a rate code rather than temporal code):

$$p(s|\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(t)}) = p(s|\bar{\mathbf{r}})
 \tag{3}$$

where $\bar{\mathbf{r}} = \frac{1}{t} \sum_{i=1}^t \mathbf{r}^{(i)}$ is the mean response after t samples (Supplemental Text). Any decoder that obeys (3) can be seen as a kind of *parametric code* over s , where the rates $\bar{\mathbf{r}}$ are the parameters. A second convenient

property for a decoder to have is if the optimal decoder is in an exponential family, or

$$p(s|\bar{\mathbf{r}}) \propto g(s) \exp(\mathbf{h}(s)^\top \mathbf{f}(\bar{\mathbf{r}})) . \quad (4)$$

Whenever (4) is true, then we would say that the neural activity forms a particular kind of parametric code called a *Nonlinear PPC* over s . A final convenient property for a decoder to have is if $\mathbf{f}(\bar{\mathbf{r}})$ is linear:

$$p(s|\bar{\mathbf{r}}) \propto g(s) \exp(\mathbf{h}(s)^\top \bar{\mathbf{r}}) . \quad (5)$$

379 This is the definition of a *Linear PPC* over s (Ma et al., 2006)¹.

380 In our simplified Encoding model, we can analytically derive the optimal decoder of the experimenter-
381 defined s conditioned on neural responses, where those neural responses are generated by sampling from the
382 Linear Gaussian internal model described above (derivation in the Supplemental Text). We find that the
383 optimal decoder is, in fact, a Linear PPC over s as defined in (5)!² This sequence of steps from equation
384 (2) through (5) suggests a general way to derive the Bayesian Decoding model implied by a given Bayesian
385 Encoding model.

386 As we discussed earlier, Encoding models typically (but not necessarily) consider inference in a task-
387 independent internal model, while Decoding models typically (but not necessarily) consider inference of
388 low-dimensional task-specific quantities. The model we described in this section is “typical” in this sense:
389 inference of \mathbf{x} is constructed to be task-independent, while the decoder of s given \mathbf{r} depends inextricably
390 on the “template function” $\mathbf{T}(s)$, which is under the control of the experimenter. The kernels $\mathbf{h}(s)$ will be
391 different for gratings of different size and spatial frequency, for plaids, or for different objects. This example
392 shows how a Bayesian Decoding model for s , implied by a task-independent Bayesian Encoding model,
393 can nonetheless be *experiment-dependent*. This points to a potentially empirically decidable question: do
394 downstream areas such as V2 interpret V1 activity like a Bayesian Decoder, or do they interpret V1 activity
395 as representing a posterior over a set of latents, \mathbf{x} ?

396 3 Discussion

397 We have identified and characterized a previously unstated difference between approaches to constructing
398 Bayesian neural models: Bayesian Encoding and Bayesian Decoding. This distinction is orthogonal to
399 existing and much-debated distinctions like whether neural responses reflect parameters or samples of the
400 inferred distribution. Making the distinction between Bayesian Encoding and Bayesian Decoding explicit
401 provides new insights into the long-standing debate about the nature of the neural code. Importantly, we
402 have demonstrated that these two approaches can give rise to different but compatible models of the same
403 neural circuit, underlining the point that Bayesian Encoding and Decoding models are complementary, and
404 not mutually exclusive. The complementary nature of these approaches has direct implications for both
405 theoretical debates and the correct interpretation of empirical data.

406 Our example model sheds light on the much-debated question of whether neural responses are more closely
407 related to parameters of the encoded probability distribution, as in probabilistic population codes (PPC; Ma
408 et al. (2006)) and in distributed distributional codes (DDC; Vertes and Sahani (2018)), or to samples from
409 the distribution as in neural sampling (reviewed in Fiser et al. (2010); Pouget et al. (2013); Sanborn (2015);
410 Gershman and Beck (2016)). In our example, the Bayesian Decoding model implies a (parametric) PPC
411 while, by construction, neural responses in the Bayesian Encoding model represent samples, demonstrating
412 that it is possible that the very same neural responses are compatible with both depending on perspective.

413 Our model is a constructive proof that Encoding and Decoding models *can* be compatible on the same
414 data, but this is will not be true in general. For instance, non-Gaussian Encoding models will not generally
415 form a linear PPC from the decoding perspective, or only for specific sets of stimuli, or they may be decodable
416 only as a nonlinear PPC. Generalizing from our specific example, the key question is, which families of
417 Encoding models, consisting of both $p_b(\mathbf{I}, \mathbf{x})$ and an assumption about the link to neural responses, are
418 compatible with which families of Decoding models, consisting of $p(s|\mathbf{r})$ and $p(\mathbf{I}|s)$? These will come in

¹PPCs also place restrictions on nuisance variables which we have omitted here.

²Further discussion of the nature of this PPC and its relation to the parameters of the internal model can be found in Shivkumar et al. (2018).

419 pairs – each family of Encoding models defines a family of compatible Decoding models, and vice versa.
420 Identifying these pairs of compatible model families is a theoretical question with important implications for
421 the interpretation of empirical data: while Encoding and Decoding models traditionally have been supported
422 (and falsified) by different kinds of empirical data (see section 2.3.4), understanding their link will allow us to
423 bridge that divide. For instance, if an Encoding model implies a particular family of Decoding models, then
424 data that falsifies that Decoding family will also falsify the Encoding family. Similarly, if a Decoding model
425 is only compatible with a family of Encoding models that is too constrained to effectively model natural
426 inputs, then that would pose a challenge for the Decoding model. As an example of this kind of argument,
427 Orbán et al. note in supplemental analyses that their sampling model appears to be empirically consistent
428 with a contrast-invariant linear PPC over orientation, but that “linear decoding of population responses will
429 significantly fall short of being optimal” once more complex tasks are considered (Orbán et al., 2016).

430 More generally, our arguments also raise questions about what makes a neural code “distributional”, i.e.
431 representing a whole distribution rather than just a point estimate, and what would constitute empirical
432 evidence for it. While the Bayesian Encoding model in our example assumed a sampling-based represen-
433 tation of the posterior over \mathbf{x} , consider a reduced, non-distributional version in which neural responses are
434 proportional to a point estimate of \mathbf{x} such as its mean or mode (Olshausen and Field, 1996). This would
435 be a poor Bayesian Encoding model in the sense that the full $p_b(\mathbf{x}|\mathbf{I})$ distribution is not recoverable from
436 \mathbf{r} . Yet, even this reduced model gives rise to a probabilistic code (PPC) over s . Such a *point estimate* code
437 over variables in the brain’s internal model would still enable many of the apparently Bayesian behaviors
438 observed in low-dimensional psychophysics tasks and used to motivate Bayesian Decoding theories, as dis-
439 cussed in section 2.3.4 above. Another example of non-Bayesian encoding but Bayesian decoding is given
440 by Orhan and Ma. It would therefore be a mistake to treat empirical evidence for near-optimal or near-
441 Bayesian behavior in a particular task alone as evidence that the brain represents probability distributions
442 over variables in an internal model of sensory inputs (Rahnev, 2019; Koblinger et al., 2021). The distinction
443 between Bayesian Encoding and Bayesian Decoding might thus productively add to the open philosophical
444 question: “if perception is probabilistic, why does it not seem probabilistic?” (Block, 2018; Rahnev et al.,
445 2020).

446 The seminal paper by Zemel et al. (1998) introduced the concept of encoding (and decoding) general
447 probability distributions in (and from) neural activity. Most work over the following 20+ years typically
448 focused on either Encoding or Decoding, as shown by Table 1, despite Zemel et al. considering both
449 perspectives as tightly linked. This divergence was likely strengthened by the fact that Encoding studies
450 almost exclusively considered internal latent variables (\mathbf{x}), while work taking the Decoding perspective
451 considered distributions over task-defined variables (s). From today’s perspective, the encoding formalism
452 of Zemel et al. and its application in Sahani and Dayan (2003) maps naturally onto our Bayesian Encoding
453 category. Furthermore, it is philosophically closely aligned with the other studies in this category, and shares
454 with them the idea that implied decoders that are non-Bayesian (but note that “decoding is only an implicit
455 operation that the system need never actually perform” Zemel et al. (1998)). Interestingly, while Zemel
456 et al. (1998) discounted the possibility of optimally decoding the encoded distribution using Bayes’ rule as
457 too inflexible, almost all later studies that took the decoding approach were based on Bayes’ rule, and now
458 form our Bayesian Decoding category.

459 The key step in our example system above which allowed us to interpret samples of \mathbf{x} as a PPC was
460 to construct the PPC over a different variable: s . This raises the question: what if s is part of the
461 brain’s internal model? One possibility is that “orientation” (or any other s in a task) is a useful abstraction
462 of natural stimuli, in which case it may have been learned (or evolved) and may permanently be a part
463 of the brain’s internal model. Another possibility is that orientation (or any other s) is part of the brain’s
464 internal model because the brain changes its internal model as the result of learning the present task (Haefner
465 et al., 2016; Lange and Haefner, 2022). Echoing section 2.3.3 above, even if s is part of the brain’s internal
466 model, Bayesian Encoding and Decoding models would nonetheless differ in their approach to the question
467 of how neural responses, \mathbf{r} , relate to the distribution on s . Bayesian Encoding models would begin with a
468 generative model of sensory input \mathbf{I} from s (and possible other internal variables \mathbf{x}) and ask how the true
469 posterior $p_b(s, \mathbf{x}|\mathbf{I})$ is represented by neural responses \mathbf{r} . Bayesian Decoding models, on the other hand,
470 would investigate the relationship between s in the world and evoked neural responses, $p(\mathbf{r}|s)$, and study a
471 different kind of posterior, $p(s|\mathbf{r})$, which takes the perspective of the experimenter, or possibly the rest of the
472 brain trying to read out s from \mathbf{r} . If the *decoded* distribution, $p(s|\mathbf{r})$, matches the ideal *encoded* distribution,

473 $p_b(s|\mathbf{I})$, then the code for s is said to be *efficient* (Beck et al., 2012).

474 The choice of variable which is assumed to be inferred, also impacts the interpretation of neural variability.
475 In our example above, neural variability is directly related to the uncertainty in the posterior over \mathbf{x} . In
476 contrast, the uncertainty over s encoded by the Bayesian decoding model is unrelated to the neural variability,
477 depending on the samples only through their *mean*, rather than their *variance*. Given sufficiently many
478 samples, the uncertainty over s is only determined by the noise in the channel between experimenter and
479 brain (Σ_{e-b}). This is an important point for experiments that seek to test the neural sampling hypothesis
480 by relating neural variability and “uncertainty”: in our example model, only uncertainty over \mathbf{x} but not over
481 s manifests as neural variability, while s is the variable most commonly and naturally manipulated in an
482 experiment.

483 The issues raised in this paper for models of visual perception also have implications for Bayesian models
484 of cognition, where ideas related to sampling (Vul and Rich, 2010; Sanborn et al., 2010; Lieder et al., 2014;
485 Vul et al., 2014; Sanborn and Chater, 2016; Lieder et al., 2017; Zhu et al., 2020), variational inference (Hohwy
486 et al., 2008; Daw et al., 2008; Sanborn and Silva, 2013), or both (Lange et al., 2021) have been invoked to
487 explain a wide variety of heuristics and biases (reviewed in Sanborn (2015); Griffiths et al. (2012b)). Here,
488 too, it is important to distinguish between probabilistic models of the world that are posited to exist in a
489 subject’s mind (as is typical in Bayesian Encoding) from experimenter-defined models of a particular task (as
490 is typical in Bayesian Decoding). Closely related is the distinction drawn by Knill and Richards between the
491 “inference problem” (what the brain infers in the internal model it assumes) and the “information problem”
492 (what information is available in the world) (Knill and Richards, 1996, Chapter 1). For example, Vul et al.
493 argue that certain deviations from Bayes-optimal behavior can be explained as the result of basing decisions
494 on a single Monte-Carlo sample. However, it is conceivable that what appears to be a single point-estimate
495 sample over a quantity relevant to a task may, in fact, be a local, perhaps unimodal distribution over a
496 detailed internal model, as in variational approximations. It is further conceivable that multiple “samples”
497 correspond to a mixture of variational approximations over an internal model (Jaakkola and Jordan, 1998;
498 Lange et al., 2022). Conversely, a single high dimensional point estimate of an internal model may be
499 sufficient to facilitate apparently Bayesian behavior with respect to a low-dimensional task. Changing our
500 reference frame from internal models to experimenter-defined tasks may make samples appear as variational
501 parameters, or vice versa.

502 Koblinger et al. (2021) recently posed the question whether uncertainty in the brain is represented “con-
503 stitutively”, i.e. about many variables regardless of their relevance for a specific task, or “opportunistically,”
504 only about task-relevant variables. While this distinction appears related to the difference between Bayesian
505 Encoding and Decoding with respect to their task-independence, there are important differences. While
506 Encoding and Decoding models have so far mostly been applied in task-independent and task-dependent
507 contexts, respectively, to what degree representations of uncertainty are task-specific is an empirical question
508 that can be productively asked within both the Encoding and Decoding approaches. For instance, Bayesian
509 Encoding models of object recognition may differ in whether they propose that the brain represents pos-
510 teriors only over task-relevant object identities, or all possible objects. One can similarly imagine both
511 “constitutive” and “opportunistic” Decoding models. For instance, the toy example model we presented
512 above is an opportunistic Decoding model, where s is determined by the experimenter, and \mathbf{r} is only said to
513 represent a distribution over s in that task context. In a constitutive Decoding model, the representation
514 of a distribution about one quantity, like $p(\text{orientation}|\mathbf{r})$, would potentially interact with representations of
515 other quantities, like $p(\text{location}|\mathbf{r})$, regardless of the immediate task-relevance of each.

516 Walker et al. (2022) pointed out a distinction between “descriptive” versus “process” approaches to the
517 study of neural representations of uncertainty. In their classification, the “descriptive” approach derives
518 estimates about the observer’s subjective uncertainty from either presented stimuli or recorded behavioral
519 reports. The “process” approach, on the other hand, derives an estimate of uncertainty from neural responses.
520 To what degree this classification is related to the Bayesian Encoding and Decoding approaches, respectively,
521 is unclear, and likely depends on additional assumptions, e.g. about the relationship between behavioral
522 reports and reference posterior in the Encoding approach, $p_b(\mathbf{x}|\mathbf{I})$, and about the nature of the model used
523 to infer uncertainty from neural responses.

524 To conclude, the Bayesian Brain Hypothesis is not a single idea, but a collection of computational
525 models, philosophical ideas, and explanations for a variety of empirical data. It is a *framework* rather than
526 a *theory* (Griffiths et al., 2012a). Bayesian Encoding and Bayesian Decoding are complementary approaches

527 to constructing concrete models within the Bayesian Brain framework. These two approaches have been a
528 major source of variation among models, and their complementary nature has previously gone unnoticed.
529 We hope that these insights will lead to clearer and more productive discussions on the nature of inference
530 in the brain, both in terms of neural representations of probability and in terms of behavior.

531 Code availability

532 Two panels in Figure 4 were generated by simulation. The code is available at <https://github.com/haefnerlab/bayesian-encoding-decoding>.
533

534 References

- 535 Laurence Aitchison and Máté Lengyel. The Hamiltonian Brain: Efficient Probabilistic Inference with
536 Excitatory-Inhibitory Neural Circuit Dynamics. *PLoS Computational Biology*, pages 1–24, 2016.
- 537 David Alais and David Burr. The Ventriloquist Effect Results from Near-Optimal Bimodal Integration.
538 *Current Biology*, 14(3):257–262, 2004.
- 539 Dora E. Angelaki, Yong Gu, and Gregory C. DeAngelis. Multisensory integration: psychophysics, neuro-
540 physiology, and computation. *Current opinion in neurobiology*, 19(4):452–8, aug 2009.
- 541 Mihály Bányai, Andreea Lazar, Liane Klein, Johanna Klon-Lipok, Marcell Stippinger, Wolf Singer, and
542 Gergő Orbán. Stimulus complexity shapes response correlations in primary visual cortex. *Proceedings of*
543 *the National Academy of Sciences*, 116(7):2723–2732, 2019.
- 544 Jeffrey M. Beck, Wei Ji Ma, Roozbeh Kiani, Timothy D. Hanks, Anne K. Churchland, Jamie Roitman,
545 Michael N. Shadlen, Peter E. Latham, and Alexandre Pouget. Probabilistic Population Codes for Bayesian
546 Decision Making. *Neuron*, 36(6):1142–1152, 2008.
- 547 Jeffrey M. Beck, Peter E. Latham, and Alexandre Pouget. Marginalization in Neural Circuits with Divisive
548 Normalization. *J. Neurosci.*, 31(43):15310–15319, 2011.
- 549 Jeffrey M. Beck, Katherine Heller, and Alexandre Pouget. Complex Inference in Neural Circuits with
550 Probabilistic Population Codes and Topic Models. *Advances in Neural Information Processing Systems*,
551 25:3068–3076, 2012.
- 552 Anthony J Bell and Terrence J Sejnowski. The "Independent Components" of Scenes are Edge Filters.
553 *Vision Research*, 37(23):3327–3338, 1997.
- 554 Pietro Berkes, Gergo Orbán, Máté Lengyel, and József Fiser. Spontaneous Cortical Activity Reveals Hall-
555 marks of an Optimal Internal Model of the Environment. *Science*, 331(January):83–87, 2011.
- 556 Johannes Bill, Lars Buesing, Stefan Habenschuss, Bernhard Nessler, Wolfgang Maass, and Robert Legenstein.
557 Distributed Bayesian computation and self-organized learning in sheets of spiking neurons with local lateral
558 inhibition. *PLoS ONE*, 10(8):1–51, 2015.
- 559 Christopher M Bishop. *Pattern Recognition and Machine Learning*. Springer, Cambridge, 2006.
- 560 Ned Block. If perception is probabilistic, why does it not seem probabilistic? *Philosophical Transactions of*
561 *the Royal Society B: Biological Sciences*, 373(1755), 2018.
- 562 Jörg Bornschein, Marc Henniges, and Jörg Lücke. Are V1 Simple Cells Optimized for Visual Occlusions? A
563 Comparative Study. *PLoS Computational Biology*, 9(6), 2013.
- 564 Lars Buesing, Johannes Bill, Bernhard Nessler, and Wolfgang Maass. Neural dynamics as sampling: A model
565 for stochastic computation in recurrent networks of spiking neurons. *PLoS Computational Biology*, 7(11),
566 2011.

- 567 Donald T. Campbell. Downward causation in Hierarchically organized biological systems. In F J Ayala,
568 editor, *Studies in the philosophy of biology*, chapter 11, pages 179–186. Macmillan Publishers Limited,
569 1974.
- 570 Ruben Coen-Cagli, Adam Kohn, and Odelia Schwartz. Flexible gating of contextual influences in natural
571 vision. *Nature Neuroscience*, 18(11):1648–1655, 2015.
- 572 Nathaniel D Daw, Aaron C Courville, and Peter Dayan. Semi-rational models of conditioning. In Nick
573 Chater and Mike Oaksford, editors, *The Probabilistic Mind:: Prospects for Bayesian cognitive science*.
574 Oxford Scholarship Online, 2008.
- 575 Peter Dayan and Larry F. Abbott. *Theoretical Neuroscience: Computational and Mathematical Modeling of*
576 *Neural Systems*. MIT Press, London, 2001.
- 577 Peter Dayan, Geoffrey E. Hinton, RM Neal, and RS Zemel. The Helmholtz machine. *Neural Computation*,
578 7(5):1–16, 1995.
- 579 Dominik Dold, Ilja Bytschok, Akos F. Kungl, Andreas Baumbach, Oliver Breitwieser, Walter Senn, Johannes
580 Schemmel, Karlheinz Meier, and Mihai A. Petrovici. Stochasticity from function — Why the Bayesian
581 brain may need no noise. *Neural Networks*, 119:200–213, 2019.
- 582 Rodrigo Echeveste, Laurence Aitchison, Guillaume Hennequin, and Máté Lengyel. Cortical-like dynamics
583 in recurrent circuits optimized for sampling-based probabilistic inference. *Nature Neuroscience*, 23:1138–
584 1149, 2020.
- 585 Marc O Ernst and Martin S Banks. Humans integrate visual and haptic information in a statistically optimal
586 fashion. *Nature*, 415(6870):429–433, 2002.
- 587 Dylan Festa, Amir Aschner, Adam Kohn, and Ruben Coen-Cagli. A Functional Model of Neuronal Response
588 Variability in Primary Visual Cortex. In *Cognitive Computational Neuroscience*, 2019.
- 589 Dylan Festa, Amir Aschner, Aida Davila, Adam Kohn, and Ruben Coen-Cagli. Neuronal variability reflects
590 probabilistic inference tuned to natural image statistics. *Nature Communications*, 12(1):1–11, 2021.
- 591 Christopher R. Fetsch, Alexandre Pouget, Gregory C. DeAngelis, and Dora E. Angelaki. Neural correlates of
592 reliability-based cue weighting during multisensory integration. *Nature Neuroscience*, 15(1):146–54, 2011.
- 593 Christopher R. Fetsch, Gregory C. DeAngelis, and Dora E. Angelaki. Bridging the gap between theories of
594 sensory cue integration and the physiology of multisensory neurons. *Nature Reviews Neuroscience*, 14(6):
595 429–442, 2013.
- 596 József Fiser, Pietro Berkes, Gergo Orbán, and Máté Lengyel. Statistically optimal perception and learning:
597 from behavior to neural representations. *Trends in cognitive sciences*, 14(3):119–30, mar 2010.
- 598 Karl J. Friston. A theory of cortical responses. *Philosophical transactions of the Royal Society of London.*
599 *Series B*, 360:815–836, 2005.
- 600 Dileep George and Jeff Hawkins. Towards a mathematical theory of cortical micro-circuits. *PLoS computa-*
601 *tional biology*, 5(10):e1000532, oct 2009.
- 602 Dileep George, Alexander Lavin, J. Swaroop Guntupalli, David Mely, Nick Hay, and Miguel Lázaro-Gredilla.
603 Cortical Microcircuits from a Generative Vision Model. In *Cognitive Computational Neuroscience*, 2018.
- 604 Samuel J. Gershman and Jeffrey M. Beck. Complex Probabilistic Inference: From Cognition to Neural
605 Computation. In Ahmed Moustafa, editor, *Computational Models of Brain and Behavior*, chapter Complex
606 Pr. Wiley-Blackwell, 2016.
- 607 Samuel J. Gershman, Edward Vul, and Joshua B. Tenenbaum. Multistability and perceptual inference.
608 *Neural Computation*, 24(1):1–24, 2012.

- 609 Agnieszka Grabska-Barwinska, Jeffrey M. Beck, Alexandre Pouget, and Peter E. Latham. Demixing odors
610 — fast inference in olfaction. *Advances in Neural Information Processing Systems*, 26, 2013.
- 611 Agnieszka Grabska-Barwińska, Simon Barthelmé, Jeff Beck, Zachary F. Mainen, Alexandre Pouget, and
612 Peter E. Latham. A probabilistic approach to demixing odors. *Nature Neuroscience*, 20(1):98–106, 2017.
- 613 Arnulf B.A. Graf, Adam Kohn, Mehrdad Jazayeri, and J. Anthony Movshon. Decoding the activity of
614 neuronal populations in macaque primary visual cortex. *Nature Neuroscience*, 14(2):239–247, 2011.
- 615 Thomas L. Griffiths, Nick Chater, Dennis Norris, and Alexandre Pouget. How the Bayesians Got Their
616 Beliefs (and What Those Beliefs Actually Are): Comment on Bowers and Davis (2012). *Psychological*
617 *Bulletin*, 138(3):415–422, 2012a.
- 618 Thomas L. Griffiths, Edward Vul, and a. N. Sanborn. Bridging Levels of Analysis for Probabilistic Models
619 of Cognition. *Current Directions in Psychological Science*, 21(4):263–268, 2012b.
- 620 Ralf M. Haefner, Pietro Berkes, and József Fiser. Perceptual Decision-Making as Probabilistic Inference by
621 Neural Sampling. *Neuron*, 90:649–660, 2016.
- 622 Guillaume Hennequin, Yashar Ahmadian, Daniel B Rubin, Máté Lengyel, and Kenneth D Miller. The Dy-
623 namical Regime of Sensory Cortex: Stable Dynamics around a Single Stimulus-Tuned Attractor Account
624 for Patterns of Noise Variability. *Neuron*, 98:846–860, 2018.
- 625 Jakob Hohwy, Andreas Roepstorff, and Karl J. Friston. Predictive coding explains binocular rivalry: An
626 epistemological review. *Cognition*, 108(3):687–701, 2008.
- 627 Han Hou, Qihao Zheng, Yuchen Zhao, Alexandre Pouget, and Yong Gu. Neural Correlates of Optimal
628 Multisensory Decision Making under Time-Varying Reliabilities with an Invariant Linear Probabilistic
629 Population Code. *Neuron*, 104:1–12, 2019.
- 630 Patrik O. Hoyer and Aapo Hyvärinen. Interpreting neural response variability as monte carlo sampling of
631 the posterior. *Advances in Neural Information Processing Systems*, 17(1):293–300, 2003.
- 632 Tommi S. Jaakkola and Michael I. Jordan. Improving the Mean Field Approximation via the Use of Mixture
633 Distributions. In Michael I. Jordan, editor, *Learning in Graphical Models*. Kluwer Academic Publishers,
634 1998.
- 635 E. T. Jaynes. *Probability Theory: The Logic of Science*. Cambridge University Press, New York, 2003.
- 636 Daniel Kersten, Pascal Mamassian, and Alan Yuille. Object perception as bayesian inference. *Annual Review*
637 *of Psychology*, pages 271–304, 2004.
- 638 David C. Knill and Alexandre Pouget. The Bayesian brain: the role of uncertainty in neural coding and
639 computation. *Trends in Neurosciences*, 27(12):712–9, dec 2004.
- 640 David C. Knill and Whitman Richards, editors. *Perception as Bayesian Inference*. Cambridge University
641 Press, New York, NY, 1996.
- 642 Ádám Koblinger, József Fiser, and Máté Lengyel. Representations of uncertainty: where art thou? *Current*
643 *Opinion in Behavioral Sciences*, 38:150–162, 2021.
- 644 Konrad P Körding. Decision Theory: What "Should" the Nervous System Do? *Science Review*, 318, 2007.
- 645 Richard D. Lange and Ralf M. Haefner. Task-induced neural covariability as a signature of approximate
646 Bayesian learning and inference. *PLoS Computational Biology*, 18(3), 2022.
- 647 Richard D. Lange, Ankani Chattoraj, Jeffrey M. Beck, Jacob L. Yates, and Ralf M. Haefner. A confirmation
648 bias in perceptual decisionmaking due to hierarchical approximate inference. *PLoS Computational Biology*,
649 17(11):1–30, 2021.

- 650 Richard D. Lange, Ari S. Benjamin, Ralf M. Haefner*, and Xaq Pitkow*. Interpolating between sampling
651 and variational inference with infinite stochastic mixtures. *UAI*, August 2022.
- 652 Alexander Lavin, J. Swaroop Guntupalli, Miguel Lázaro-gredilla, Wolfgang Lehrach, and Dileep George.
653 Explaining Visual Cortex Phenomena using Recursive Cortical Network. In *Cognitive Computational*
654 *Neuroscience*, 2018.
- 655 Falk Lieder, Ming Hsu, and Thomas L. Griffiths. The high availability of extreme events serves resource-
656 rational decision-making. In *Cognitive Science Society*, pages 2567–2572, 2014.
- 657 Falk Lieder, Thomas L. Griffiths, Quentin J M Huys, and Noah D. Goodman. The anchoring bias reflects
658 rational use of cognitive resources. *Psychonomic Bulletin & Review*, 2017.
- 659 Wei Ji Ma, Jeffrey M. Beck, Peter E. Latham, and Alexandre Pouget. Bayesian inference with probabilistic
660 population codes. *Nature Neuroscience*, 9(11):1432–1438, 2006.
- 661 R. Moreno-Bote, David C. Knill, and A. Pouget. Bayesian sampling in visual perception. *Proceedings of the*
662 *National Academy of Sciences*, 108(30):12491–12496, 2011.
- 663 Richard Morris. Developments of a water-maze procedure for studying spatial learning in the rat. *Journal*
664 *of Neuroscience Methods*, 11(1):47–60, 1984.
- 665 Kevin P. Murphy. *Machine Learning: A Probabilistic Perspective*. The MIT Press, Cambridge, MA, 2012.
- 666 Bruno A. Olshausen and David J. Field. Emergence of simple-cell receptive field properties by learning a
667 sparse code for natural images, 1996.
- 668 Bruno a Olshausen and David J. Field. Sparse coding with an incomplete basis set: a strategy employed by
669 V1?, 1997.
- 670 Gergő Orbán, Pietro Berkes, József Fiser, and Máté Lengyel. Neural Variability and Sampling-Based Prob-
671 abilistic Representations in the Visual Cortex. *Neuron*, 92(2):530–543, 2016.
- 672 A. Emin Orhan and Wei Ji Ma. Efficient probabilistic inference in generic neural networks trained with
673 non-probabilistic feedback. *Nature Communications*, 8(138):1–30, 2017.
- 674 Dejan Pecevski, Lars Buesing, and Wolfgang Maass. Probabilistic inferences general graphical models
675 through sampling in stochastic networks of spiking neurons. *PLOS Computational Biology*, 7(12), 2011.
- 676 Mihai A. Petrovici, Johannes Bill, Ilja Bytschok, Johannes Schemmel, and Karlheinz Meier. Stochastic
677 inference with spiking neurons in the high-conductance state. *Physical Review E*, 94, 2016.
- 678 Alexandre Pouget, Jeffrey M. Beck, Wei Ji Ma, and Peter E. Latham. Probabilistic brains: knowns and
679 unknowns. *Nature Neuroscience*, 16(9):1170–8, 2013.
- 680 Dimitri Probst, Mihai A. Petrovici, Ilja Bytschok, Johannes Bill, Dejan Pecevski, Johannes Schemmel, and
681 Karlheinz Meier. Probabilistic inference in discrete spaces can be implemented into networks of LIF
682 neurons. *Frontiers in computational neuroscience*, 9(13):1–11, 2015.
- 683 Dobromir Rahnev. The bayesian brain: What is it and do humans have it? *Behavioral and Brain Sciences*,
684 42:e238, 2019.
- 685 Dobromir Rahnev, Ned Block, Janneke Jehee, and Rachel Denison. Is perception probabilistic? In *Cognitive*
686 *Computational Neuroscience*, 2020.
- 687 Rajkumar V. Raju and Xaq Pitkow. Inference by Reparameterization in Neural Population Codes. *Advances*
688 *in Neural Information Processing Systems*, 30, 2016.
- 689 Maneesh Sahani and Peter Dayan. Doubly Distributional Population Codes: Simultaneous Representation
690 of Uncertainty and Multiplicity. *Neural Computation*, 15:2255–2279, 2003.

- 691 Adam N Sanborn. Types of approximation for probabilistic cognition: Sampling and variational. *Brain and*
692 *Cognition*, 2015.
- 693 Adam N Sanborn and Nick Chater. Bayesian Brains without Probabilities. *Trends in Cognitive Sciences*, 20
694 (12):883–893, 2016.
- 695 Adam N. Sanborn and Ricardo Silva. Constraining bridges between levels of analysis: A computational
696 justification for locally Bayesian learning. *Journal of Mathematical Psychology*, 57(3-4):94–106, 2013.
- 697 Adam N Sanborn, Thomas L Griffiths, and Daniel J Navarro. Rational approximations to rational models:
698 alternative algorithms for category learning. *Psychological Review*, 117(4):1144–67, oct 2010.
- 699 Cristina Savin and Sophie Denève. Spatio-temporal representations of uncertainty in spiking neural networks.
700 *Advances in Neural Information Processing Systems*, 2014.
- 701 Odelia Schwartz and Eero P Simoncelli. Natural signal statistics and sensory gain control. *Nature Neuro-*
702 *science*, 4(8):819–825, 2001.
- 703 Sabyasachi Shivkumar, Richard D. Lange, Ankani Chattoraj, and Ralf M. Haefner. A probabilistic population
704 code based on neural samples. *Advances in Neural Information Processing Systems*, 31:7070–7079, 2018.
- 705 Chihiro I. Tajima, Satohiro Tajima, Kowa Koida, Hidehiko Komatsu, Kazuyuki Aihara, and Hideyuki Suzuki.
706 Population code dynamics in categorical perception. *Nature Scientific Reports*, 6(22536):1–13, 2016.
- 707 Eszter Vertes and Maneesh Sahani. Flexible and accurate inference and learning for deep generative models.
708 *Advances in Neural Information Processing Systems*, 31, 2018.
- 709 Hermann von Helmholtz. *Treatise on Physiological Optics*. The Optical Society of America, 1925.
- 710 Edward Vul and Anina N. Rich. Independent Sampling of Features Enables Conscious Perception of Bound
711 Objects. *Psychological Science*, 21(8):1168–1175, 2010.
- 712 Edward Vul, Noah D. Goodman, Thomas L. Griffiths, and Joshua B. Tenenbaum. One and done? Optimal
713 decisions from very few samples. *Cognitive Science*, 38(4):599–637, 2014.
- 714 Martin J. Wainwright and Michael I. Jordan. Graphical Models, Exponential Families, and Variational
715 Inference. *Foundations and Trends in Machine Learning*, 1(1–2):1–305, 2008.
- 716 Edgar Y Walker, R. James Cotton, Wei Ji Ma, and Andreas S Tolias. A neural basis of probabilistic
717 computation in visual cortex. *Nature Neuroscience*, 23:122–129, 2019.
- 718 Edgar Y Walker, Stephan Pohl, Rachel N Denison, David L Barack, Jennifer Lee, Ned Block, Wei Ji Ma,
719 and Florent Meyniel. Studying the neural representations of uncertainty. *arXiv*, pages 1–28, 2022.
- 720 Stephen Yablo. Mental Causation. *The Philosophical Review*, 101(2):245–280, 1992.
- 721 Richard S. Zemel, Peter Dayan, and Alexandre Pouget. Probabilistic Interpretation of Population Codes.
722 *Neural Computation*, 10(2):403–430, 1998.
- 723 Jian-Qiao Zhu, Adam N Sanborn, and Nick Chater. The Bayesian Sampler: Generic Bayesian Inference
724 Causes Incoherence in Human Probability Judgments. *Psychological Review*, 127(5):719–748, 2020.

725 S Supplemental Text

726 S.1 Derivation of decoded posterior for Gaussian prior over \mathbf{x}

727 In this section we provide a brief derivation of the optimal posterior over an experimenter-defined s con-
 728 ditioned on samples of internal-model variables \mathbf{x} , where the brain’s internal model $p_b(\mathbf{x}, \mathbf{I})$ is assumed to
 729 be a linear Gaussian model with a Gaussian prior over \mathbf{x} . The use of a Gaussian prior over \mathbf{x} is a further
 730 simplification of the derivation in Shivkumar et al. (2018). Formally, the setup is as follows:

- 731 1. Assume that the scalar s (such as orientation) gives rise to observed images \mathbf{I} as

$$\mathbf{I} = \mathbf{T}(s) + \boldsymbol{\eta},$$

732 where $\mathbf{T}(s)$ is a “template” function (such as a grating image), and $\boldsymbol{\eta}$ is zero-mean Gaussian-distributed
 733 pixel noise with covariance $\boldsymbol{\Sigma}_{e-b}$.

- 734 2. Assume that the brain’s internal model, $p_b(\mathbf{x}, \mathbf{I})$, factorizes as $p_b(\mathbf{x})p_b(\mathbf{I}|\mathbf{x})$, where the prior is Gaussian,

$$p_b(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p),$$

735 and images are assumed to be generated as a linear combination of basis vectors,

$$p_b(\mathbf{I}|\mathbf{x}) = \mathcal{N}(\mathbf{I}; \mathbf{A}\mathbf{x}, \boldsymbol{\Sigma}_\mathbf{I}).$$

- 736 3. **Bayesian Encoding Model:** Assume that, conditioned on \mathbf{I} , the brain samples $\{\mathbf{x}^{(i)}\} \sim p_b(\mathbf{x}|\mathbf{I})$,
 737 and that each value of \mathbf{x} corresponds to a neuron, so that $\mathbf{r}^{(i)} = \mathbf{x}^{(i)}$. (We will use “ $\mathbf{r}^{(i)}$ ” to denote the
 738 vector of neural activity at time i , and $\mathbf{r} = \{\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(t)}\}$ to denote all neural activity in the relevant
 739 population up to time t .)

- 740 4. **Bayesian Decoding Model:** We will derive the Bayesian decoder of s given $\{\mathbf{r}^{(1)}, \dots, \mathbf{r}^{(t)}\}$.

741 We are interested in the optimal decoder of s after t time has elapsed, or $p(s|\mathbf{r})$. By Bayes’ rule, this
 742 is proportional to $p(\mathbf{r}|s)p(s) = p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}|s)p(s)$. That is, the quantity we must compute in order to
 743 optimally decode $p(s|\mathbf{r})$ is the probability of seeing a given set of samples, $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$, provided a value of
 744 s .

745 Since s affects \mathbf{r} through \mathbf{I} , this likelihood function can be evaluated by marginalizing across all possible
 746 images

$$p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}|s) = \int p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}|\mathbf{I})p(\mathbf{I}|s)d\mathbf{I} \quad (\text{S1})$$

747 We know from our definition that

$$p(\mathbf{I}|s) = \mathcal{N}(\mathbf{I}; \mathbf{T}(s), \boldsymbol{\Sigma}_{e-b}),$$

748 and the posterior probability of all t independent samples for a given \mathbf{I} is

$$p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}|\mathbf{I}) = \prod_{i=1}^t p(\mathbf{x}^{(i)}|\mathbf{I}).$$

749 Under the simplifying assumption that both $p_b(\mathbf{x})$ and $p_b(\mathbf{I}|\mathbf{x})$ are Gaussian, the brain’s internal model
 750 posterior, $p_b(\mathbf{x}|\mathbf{I})$ is also Gaussian,

$$p(\mathbf{x}^{(i)}|\mathbf{I}) = \mathcal{N}(\mathbf{x}^{(i)}; \boldsymbol{\mu}', \boldsymbol{\Sigma}'), \quad (\text{S2})$$

where

$$\begin{aligned} \boldsymbol{\mu}' &= \boldsymbol{\Sigma}'(\boldsymbol{\Sigma}_p^{-1}\boldsymbol{\mu}_p + \mathbf{A}^\top\boldsymbol{\Sigma}_\mathbf{x}^{-1}\mathbf{I}) && \text{and} \\ \boldsymbol{\Sigma}' &= (\boldsymbol{\Sigma}_p^{-1} + \mathbf{A}^\top\boldsymbol{\Sigma}_\mathbf{x}^{-1}\mathbf{A})^{-1}. \end{aligned}$$

751 Note that the only dependence on \mathbf{I} (and therefore on s) is through $\boldsymbol{\mu}'$.

Equation (S2) gives the probability of seeing a single sample $\mathbf{x}^{(i)}$ given \mathbf{I} . The probability of all t samples is

$$\begin{aligned} p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)} | \mathbf{I}) &= \prod_{i=1}^t \mathcal{N}(\mathbf{x}^{(i)}; \boldsymbol{\mu}', \boldsymbol{\Sigma}') \\ &= \mathcal{N}(\bar{\mathbf{x}}; \boldsymbol{\mu}', t^{-1} \boldsymbol{\Sigma}') c(\mathbf{x}, \boldsymbol{\Sigma}') \end{aligned}$$

752 where $\bar{\mathbf{x}} = t^{-1} \sum_{i=1}^t \mathbf{x}^{(i)}$ is the average of samples up to time t , and $c(\mathbf{x}, \boldsymbol{\Sigma}')$ is a term that depends on
753 $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$ and on $\boldsymbol{\Sigma}'$ but not on $\boldsymbol{\mu}'$ (and therefore not on \mathbf{I} , so it can be dropped later).

754 We can now evaluate the integral in (S1) to get the probability of t samples for a given s :

$$\begin{aligned} 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} \\ &= \int \mathcal{N}(\bar{\mathbf{x}}; \boldsymbol{\mu}', t^{-1} \boldsymbol{\Sigma}') c(\mathbf{x}, \boldsymbol{\Sigma}') \mathcal{N}(\mathbf{I}; \mathbf{T}(s), \boldsymbol{\Sigma}_{e-b}) d\mathbf{I} \\ &= c(\mathbf{x}, \boldsymbol{\Sigma}') \int \mathcal{N}(\bar{\mathbf{x}}; \underbrace{\boldsymbol{\Sigma}'(\boldsymbol{\Sigma}_p^{-1} \boldsymbol{\mu}_p + \mathbf{A}^\top \boldsymbol{\Sigma}_x^{-1} \mathbf{I})}_{\text{definition of } \boldsymbol{\mu}'}, t^{-1} \boldsymbol{\Sigma}') \mathcal{N}(\mathbf{I}; \mathbf{T}(s), \boldsymbol{\Sigma}_{e-b}) d\mathbf{I} \\ &= c(\mathbf{x}, \boldsymbol{\Sigma}') \int \mathcal{N}(\bar{\mathbf{x}} - \boldsymbol{\Sigma}' \boldsymbol{\Sigma}_p^{-1} \boldsymbol{\mu}_p; \underbrace{\boldsymbol{\Sigma}' \mathbf{A}^\top \boldsymbol{\Sigma}_x^{-1} \mathbf{I}}_{\text{Let } \mathbf{x}' \equiv \boldsymbol{\Sigma}' \mathbf{A}^\top \boldsymbol{\Sigma}_x^{-1} \mathbf{I}}, t^{-1} \boldsymbol{\Sigma}') \mathcal{N}(\mathbf{I}; \mathbf{T}(s), \boldsymbol{\Sigma}_{e-b}) d\mathbf{I} \\ &(*) \propto \int \mathcal{N}(\bar{\mathbf{x}} - \boldsymbol{\Sigma}' \boldsymbol{\Sigma}_p^{-1} \boldsymbol{\mu}_p; \mathbf{x}', t^{-1} \boldsymbol{\Sigma}') \mathcal{N}(\mathbf{x}'; \boldsymbol{\Sigma}' \mathbf{A}^\top \boldsymbol{\Sigma}_x^{-1} \mathbf{T}(s), \boldsymbol{\Sigma}' \mathbf{A}^\top \boldsymbol{\Sigma}_x^{-1} \boldsymbol{\Sigma}_{e-b} \boldsymbol{\Sigma}_x^{-1} \mathbf{A} \boldsymbol{\Sigma}') d\mathbf{x}' \\ &= \mathcal{N}(\bar{\mathbf{x}} - \boldsymbol{\Sigma}' \boldsymbol{\Sigma}_p^{-1} \boldsymbol{\mu}_p; \boldsymbol{\Sigma}' \mathbf{A}^\top \boldsymbol{\Sigma}_x^{-1} \mathbf{T}(s), t^{-1} \boldsymbol{\Sigma}' + \boldsymbol{\Sigma}' \mathbf{A}^\top \boldsymbol{\Sigma}_x^{-1} \boldsymbol{\Sigma}_{e-b} \boldsymbol{\Sigma}_x^{-1} \mathbf{A} \boldsymbol{\Sigma}') \\ &= \mathcal{N}(\bar{\mathbf{x}}; \boldsymbol{\mu}'', \boldsymbol{\Sigma}'') \end{aligned}$$

where

$$\begin{aligned} \boldsymbol{\mu}'' &= \boldsymbol{\Sigma}' (\boldsymbol{\Sigma}_p^{-1} \boldsymbol{\mu}_p + \mathbf{A}^\top \boldsymbol{\Sigma}_x^{-1} \mathbf{T}(s)) \\ \boldsymbol{\Sigma}'' &= t^{-1} \boldsymbol{\Sigma}' + \boldsymbol{\Sigma}' \mathbf{A}^\top \boldsymbol{\Sigma}_x^{-1} \boldsymbol{\Sigma}_{e-b} \boldsymbol{\Sigma}_x^{-1} \mathbf{A} \boldsymbol{\Sigma}' \end{aligned}$$

755 In the line marked (*), we changed variables to switch from an integral over \mathbf{I} to an integral over \mathbf{x}' . This
756 line is a proportionality because we also dropped terms that do not depend on s , including the Jacobian
757 term from the change of variables, since later we will use this expression as a likelihood function of s .

Expanding the definition of $\mathcal{N}(\dots)$, we can now write the posterior over s given $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$ as

$$\begin{aligned} p(s | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}) &= p(s | \bar{\mathbf{x}}) \\ &\propto p(s) \exp\left(-\frac{1}{2} (\bar{\mathbf{x}} - \boldsymbol{\mu}'')^\top \boldsymbol{\Sigma}''^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}'')\right) \\ &\propto p(s) \exp\left(\bar{\mathbf{x}}^\top \boldsymbol{\Sigma}''^{-1} \boldsymbol{\mu}'' - \frac{1}{2} \boldsymbol{\mu}''^\top \boldsymbol{\Sigma}''^{-1} \boldsymbol{\mu}''\right). \end{aligned}$$

Substituting $\bar{\mathbf{r}}$ for $\bar{\mathbf{x}}$ and rewriting in terms of a Linear PPC, this is

$$p(s | \mathbf{r}) \propto g(s) \exp(\mathbf{h}(s)^\top \bar{\mathbf{r}}) \quad (5) \text{ restated}$$

where

$$\begin{aligned} g(s) &= p(s) \exp\left(-\frac{1}{2} \boldsymbol{\mu}''(s)^\top \boldsymbol{\Sigma}''^{-1} \boldsymbol{\mu}''(s)\right) \\ \mathbf{h}(s) &= \boldsymbol{\Sigma}''^{-1} \boldsymbol{\Sigma}' \mathbf{A}^\top \boldsymbol{\Sigma}_x^{-1} \mathbf{T}(s). \end{aligned}$$