

Bayesian sampling in visual perception

Rubén Moreno-Bote^{a,b,1}, David C. Knill^{b,c}, and Alexandre Pouget^b

^aFoundation Sant Joan de Déu, Parc Sanitari Sant Joan de Déu, Esplugues de Llobregat, 08950 Barcelona, Spain; and ^bDepartment of Brain and Cognitive Sciences and ^cCenter for Visual Science, University of Rochester, Rochester, NY, 14627

Edited by Wilson S. Geisler, University of Texas at Austin, Austin, TX, and approved May 31, 2011 (received for review January 27, 2011)

It is well-established that some aspects of perception and action can be understood as probabilistic inferences over underlying probability distributions. In some situations, it would be advantageous for the nervous system to sample interpretations from a probability distribution rather than commit to a particular interpretation. In this study, we asked whether visual percepts correspond to samples from the probability distribution over image interpretations, a form of sampling that we refer to as Bayesian sampling. To test this idea, we manipulated pairs of sensory cues in a bistable display consisting of two superimposed moving drifting gratings, and we asked subjects to report their perceived changes in depth ordering. We report that the fractions of dominance of each percept follow the multiplicative rule predicted by Bayesian sampling. Furthermore, we show that attractor neural networks can sample probability distributions if input currents add linearly and encode probability distributions with probabilistic population codes.

Bayesian inference | neuronal network | neuronal noise | perceptual bistability

There is mounting evidence that neural circuits can implement probabilistic inferences over sensory, cognitive, or motor variables. In some cases, humans can perform these inferences optimally, as in multi-cue or multisensory integration (1–8). For complex tasks, such as object recognition, action perception, and object tracking, the computations required for optimal inference are intractable, which implies that humans must use approximate inferences (9–11). One approximate scheme that is particularly appealing from a biological point of view is sampling. Consider as an example the problem of object recognition. The goal of the inference in this case would be to compute the probability over object identities given the image. Although this probability distribution may be difficult to compute explicitly, one can often design algorithms to generate samples from the distribution, allowing one to perform approximate inference (12, 13). Some human cognitive choice behaviors suggest that the nervous system implements sampling. However, whether the same is true for low-level perceptual processing is currently unknown.

Stimuli that lead to bistable percepts (14–18), like the Necker cube, provide a tractable experimental preparation for testing the sampling hypothesis. With such stimuli, perception alternates stochastically between two possible interpretations, a behavior consistent with sampling as suggested by several works (16, 19, 20). However, the key question is what probability distribution is being sampled. If the brain uses sampling for Bayesian inference, neural circuits should sample from an internal probability distribution on possible stimulus interpretations that are conditioned on the available sensory data, the so-called posterior distribution. This distribution places important constraints on the distributions of perceptual states for bistable stimuli.

To test this idea, we used stimuli composed of two drifting gratings whose depth ordering is ambiguous (21). We then manipulated two depth cues to vary the fractions of dominance of the percepts. Our central prediction is that the fractions of dominance of each percept should behave as probabilities if they are the result of a sampling process of a posterior distribution over image interpretations. We will refer to this form of sampling as Bayesian sampling. First, we show that subjects' fractions of dominance in different cue conditions follow the same multiplicative rule as

probabilities in the Bayesian calculus, suggesting that bistable perception is indeed a form of Bayesian sampling. Second, we describe possible neural implementations of a Bayesian sampling process using attractor networks, and we discuss the link with probabilistic population codes (22).

Results

Multiplicative Rule for Combining Empirical Fractions of Dominance.

We asked subjects to report their spontaneous alternations in perceived depth ordering of two superimposed moving gratings over a 1-min period and measured the fraction of dominance time for each percept (*Methods* and Fig. 1*A*). In the first experiment, the two drifting gratings, α and β , were parameterized by their wavelength and speed. One of the wavelengths was always set to a fixed value λ^* , and one of the speeds was set to a fixed value ν^* . The remaining wavelength and speed parameters, λ and ν , respectively, determined the difference in wavelength and speed between gratings α and β , denoted $\Delta\lambda$ and $\Delta\nu$, and hence, the information for choosing grating α as the one behind. We refer to these differences as the cues to depth ordering, and we refer to the condition where the two differences are zero as the neutral cue condition ($\Delta\lambda = 0$ and $\Delta\nu = 0$). These cues have been shown to have a strong effect on the depth ordering of the gratings because of their relationship with the natural statistics of wavelength and speed of distant objects (21). In the second experiment, we manipulated wavelength and disparity, d , of the gratings. In this case, the label ν should be interchanged with the label d .

According to the Bayesian sampling hypothesis, the empirical fractions of dominance arise from a process that samples the posterior distribution on possible scene interpretations given the sensory input. As we show in *SI Methods*, when two conditionally independent cues are available (i.e., the values of the cues are independent when conditioned on true depth), an optimal system should sample from a probability distribution given by the normalized product of the probability distributions derived by varying each cue in isolation while keeping the other cue neutral. Our hypothesis implies that the empirical fractions should behave as probabilities, and therefore, they should follow the multiplicative rule (Eq. 1)

$$f_{\lambda\nu} = \frac{f_{\lambda}f_{\nu}}{f_{\lambda}f_{\nu} + (1-f_{\lambda})(1-f_{\nu})}, \quad [1]$$

where $f_{\lambda\nu}$ is the fraction of time that subjects report percept A (grating α moving behind grating β) when the cues are set to $\Delta\lambda$ and $\Delta\nu$, f_{λ} is the fraction of dominance of percept A when the speed cue is neutral ($\Delta\nu = 0$) while the wavelength cue has value $\Delta\lambda$, and f_{ν} is the dominance fraction when the wavelength cue is neutral ($\Delta\lambda = 0$) while the speed cue has value $\Delta\nu$. This relation holds whether subjects are sampling from posterior distributions

Author contributions: R.M.-B., D.C.K., and A.P. designed research; R.M.-B., D.C.K., and A.P. performed research; R.M.-B. analyzed data; and R.M.-B., D.C.K., and A.P. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

¹To whom correspondence should be addressed. E-mail: rmoreno@bcs.rochester.edu.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1101430108/-DCSupplemental.

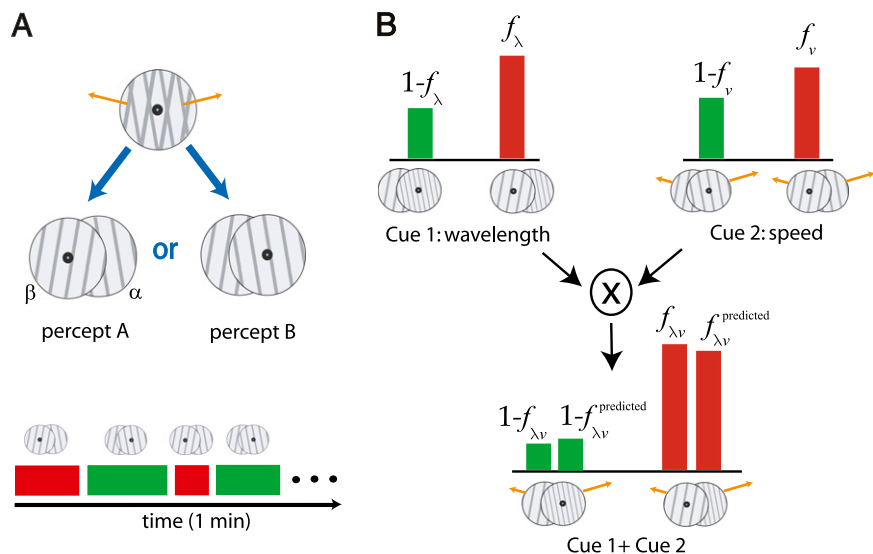


Fig. 1. Cue combination in a perceptually bistable stimulus. (A) The visual stimulus consisted of two superimposed drifting gratings moving in different directions. The perceived depth ordering of the gratings is bistable. We measured the fraction of dominance of each percept by asking subjects to report the perceived depth ordering of the gratings during trials of 1-min duration (hypothetical trial shown). (B) Cue combination. (Upper Left) Fractions of dominance for each depth ordering when wavelength is nonneutral (its value differs between the two gratings), whereas speed is neutral (its value is identical across gratings). (Upper Right) is the same as Upper Left, but when wavelength is neutral, the speed is nonneutral. (Lower) Fraction of dominance when both speed and wavelength are nonneutral. Bayesian sampling predicts that the fraction of dominance when both cues are nonneutral is equal to the normalized product of the fractions of dominance when only one cue is nonneutral (Eq. 1). In the example illustrated here, both cues were congruent.

on depth or posterior distributions raised to an arbitrary power n (*SI Methods*). The multiplicative rule provides an empirical consistency constraint for Bayesian sampling. Note that this rule does not specify how the samples are extracted over time [i.e., it works whether the samples are independent over time (23, 24) or correlated]. As discussed later, bistable perception is only consistent with a sampling mechanism that generates correlated samples (i.e., the percept tends to remain the same over hundreds of milliseconds).

Observed vs. Predicted Fractions of Dominance. The multiplicative rule was tested in two experiments. In the first experiment, the wavelength and speed differences between the two gratings, $\Delta\lambda$ and Δv , were changed from trial to trial congruently [C condition (i.e., both cues favoring the same depth ordering); example in Fig. 1B] or incongruently [IC condition (i.e., the cues favored different depth orderings)]. This change was achieved by decreasing the wavelength and increasing the speed of grating α in the C condition, while decreasing the wavelength of grating α and increasing the speed of grating β in the IC condition. In the second experiment, the wavelength and stereo disparity (instead of speed) of the gratings were manipulated in the C and IC conditions as in the previous experiment.

As shown in Figs. 2 and 3, wavelength, speed, and disparity differences in the gratings have a strong impact on the fractions of dominance of the gratings' depth ordering (21). The fraction of dominance of percept A (grating α is behind grating β) increases as the wavelength difference between gratings α and β ($\Delta\lambda = \lambda_\alpha - \lambda_\beta$) decreases. The fraction increases as the speed difference between gratings α and β ($\Delta v = v_\alpha - v_\beta$) increases in the C condition (Fig. 2A). Conversely, the fraction decreases as the difference (in speed or wavelength) between the gratings decreases in the IC condition (Fig. 2B). In the second experiment, the fraction of dominance of percept A increases as the disparity difference between gratings α and β ($\Delta d = d_\alpha - d_\beta$) increases in the C condition (Fig. 3A). Again, the reverse pattern is observed in the IC condition (Fig. 3B). In the two experiments, when the two cues are set to their neutral values, the fractions (Figs. 2 and 3, black open circles) are not significantly different from one-half [two-tailed t test; experiment 1: $p = 0.39$ (C), $p = 0.06$ (IC) and experiment 2: $p = 0.31$ (C), $p = 0.051$ (IC)].

The experimental results were compared with the theoretical predictions from the multiplicative rule (Eq. 1) (Figs. 2A and B and 3A and B). The predictions when the two cues are nonneutral (Figs. 2A and B and 3A and B, filled blue circles) were

computed using the experimental data of the single nonneutral cue cases only (Figs. 2A and B and 3A and B, open red circles). The case in which wavelength is the only nonneutral cue corresponds to the lower line of open circles in Figs. 2A and 3A and the upper line in Figs. 2B and 3B in both experiments. The cases in which speed (or disparity) is the only nonneutral cue correspond to the vertical line of open circles in the wavelength and speed (or disparity) experiment in Fig. 2B (Fig. 3B respectively). The match between the observed data points (filled red circles) and predictions is tight, even though the multiplicative rule is parameter-free and cannot be adjusted to match the experimental results (note that, for the sake of clarity, the blue dots have been slightly displaced to the right). The data in Figs. 2A and B and 3A and B were replotted in Figs. 2C and 3C to show the predicted fraction of dominance from the multiplicative model vs. the observed fraction when the two cues were nonneutral with the C (Figs. 2C and 3C, light blue dots) and IC (Figs. 2C and 3C, dark blue) conditions combined. The strong alignment of the data points along the unity line confirms that the multiplicative rule provides a tight fit to the data. Individual subjects also followed the multiplicative rule (*SI Methods* and Fig. S1).

We also tested alternative models to the multiplicative rule. In the first model, we assumed that integration between the cues does not take place—a strongest cue take all model. In this model, performance is driven by the cue with the lowest uncertainty: The fraction of dominance when both cues are varied together is set to that of the cue whose fraction when the cues are manipulated alone has the largest absolute value difference with respect to one-half (*SI Methods*). As shown in Figs. 2D and 3D (brown dots) this model fails to capture our experimental results. In the second model, we generated predictions from a realistic neuronal network (see *Results, Sampling with Realistic Neural Circuits*). When the input neurons to the network fired nonlinearly in response to the stimuli (25), the predictions of the model, which fit the single nonneutral cue conditions, substantially differed from the experimental data in the four nonneutral cues conditions (NL net) (Figs. 2D and 3D, orange dots). When the input neurons fired linearly (26), the predictions were identical to the multiplicative rule (L net) (Figs. 2D and 3D, blue dots). This result shows that the mere fact that a network can oscillate stochastically between two percepts in a way suggestive of sampling does not guarantee that it will also follow the multiplicative rule. Whether it does depends critically on how the inputs are combined, a point that we discuss more thoroughly below.

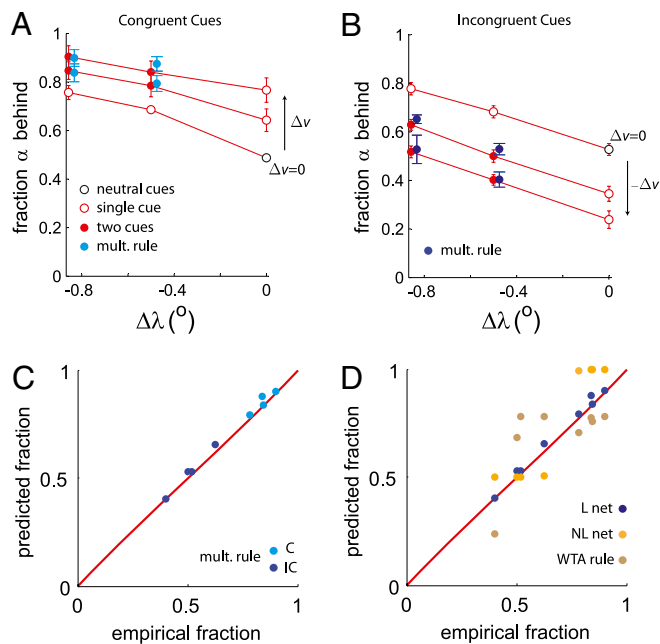


Fig. 2. Experimental and predicted fractions of dominance in the wavelength and speed cue combination experiment. (A) Fraction of dominance of percept A (i.e., grating α is behind grating β) as a function of the wavelength difference between gratings α and β ($\Delta\lambda = \lambda_\alpha - \lambda_\beta$) for three different speed differences ($\Delta v = v_\alpha - v_\beta$) in the congruent condition (both cues favored the same depth ordering). Data are averaged across subjects, and the error bars correspond to SEM across subjects. Experimental observations (red and black) and predictions from the multiplicative rule (blue circles) (Eq. 1) are shown. The predictions from the multiplicative rule were computed using the experimental data from the conditions in which only one cue was nonneutral (open circles). The black open circles correspond to the fractions measured when the two cues were neutral. The predictions are displaced slightly right in relation to the experimental data (filled red circles) to allow better visual comparison. (B) Same as in A but for the incongruent condition (the cues favored opposite percepts). (C) Predicted fractions of dominance for the multiplicative rule combining the data from the congruent (C; light blue) and incongruent (IC; dark blue) conditions from A and B as a function of the empirical fractions. (D) Same as in C but for the strongest cue take all rule (brown) and a rate-based model with nonlinear (orange) and linear (blue) input neurons.

Diffusion in an Energy Model. Our finding that bistable perception behaves like a Bayesian sampling process raises the issue as to how neurons could implement such a process. We first show that implementing the multiplicative rule is surprisingly straightforward with energy models. In *Results, Sampling with Realistic Neural Circuits*, we will present a neural instantiation of this conceptual framework. We model the dynamics of two neural populations, A and B , whose states are described by their firing rates r_A and r_B , respectively (Fig. 4A). The reduced dynamics tracks the difference between the firing rates, $r = r_A - r_B$, where $r > 0$ corresponds to percept A . This variable obeys (Eq. 2)

$$\tau \frac{d}{dt} r = -4r(r^2 - 1) + g(I_\lambda, I_v) + n(t), \quad [2]$$

where $g(I_\lambda, I_v)$ is a bias provided by the inputs and $n(t)$ is a filtered white noise with variance σ^2 (27) (*SI Methods*). The first term on the right-hand side ensures that the activity difference, r , hovers around the centers of the two energy wells (Fig. 4B). The bias term measures the combined strength of the cues, which is a function of the individual strengths I_λ and I_v , favoring percept A from the wavelength and speed cues, respectively. The function $g(I_\lambda, I_v)$ is chosen such that it is zero when the two cues are neutral (zero

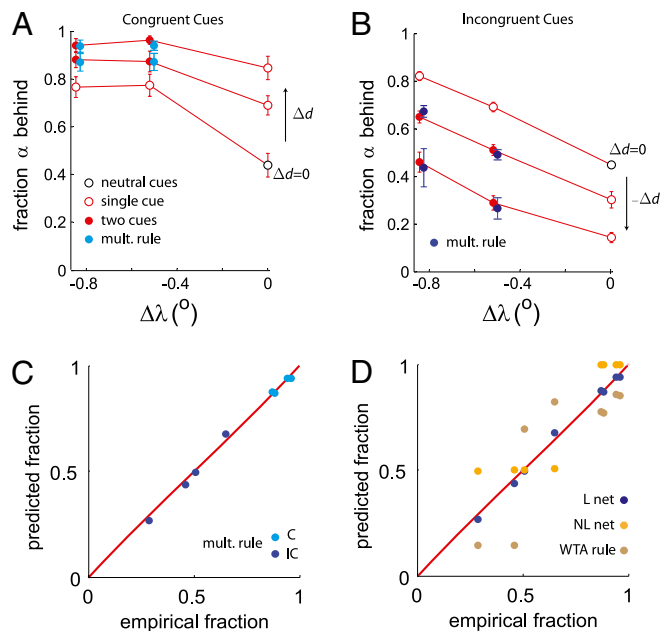


Fig. 3. Experimental and predicted fractions of dominance in the wavelength and disparity cue combination experiment. A–D are the same as in Fig. 2 but with speed replaced by disparity.

currents) and positive when the two cues favor percept A (the two currents are positive). The dynamics of Eq. 2 can be viewed as a noisy descent over the energy landscape $E(r) = r^2(r^2 - 2) - g(I_\lambda, I_v)r$, which is symmetrical (Fig. 4B, black line) when the two cues are neutral and negatively tilted (Fig. 4B, gray line) when the cues favor percept A . The resulting dynamics effectively draws samples from an underlying probability distribution that depends on the input currents (a process known as Langevin Monte Carlo sampling) (28).

To model the experimental data that we have described, we need a form of sampling that obeys the multiplicative rule. Whether the network obeys the rule or not depends critically on the function $g(I_\lambda, I_v)$. We consider here the family of functions described by $g(I_\lambda, I_v) = I_\lambda + I_v + \varepsilon(I_\lambda^2 I_v + I_v^2 I_\lambda)$, where ε measures the strength of the nonlinearity. Similar nonlinear functional dependences on the input currents naturally arise in neuronal networks with nonlinear activation functions (*Results, Sampling with Realistic Neural Circuits*).

For a value of ε different from zero, the dynamical system does not follow the multiplicative rule (Fig. 4D). In contrast, if we set ε to zero, such that $g(I_\lambda, I_v) = I_\lambda + I_v$, the system now obeys the multiplicative rule (Fig. 4E). This result can be derived analytically by computing the mean dominance duration of each percept, which corresponds to the mean escape time from one of the energy wells (*SI Methods*). We can then show that the fraction of dominance of population A for ε equal to zero is a sigmoid function of the sum of the inputs (Eq. 3)

$$f_{\lambda v} = f(s = A | I_\lambda, I_v) = \frac{1}{1 + e^{-2(I_\lambda + I_v)/\sigma_{\text{eff}}^2}} \propto e^{(I_\lambda + I_v)/\sigma_{\text{eff}}^2}, \quad [3]$$

where σ_{eff}^2 is the effective noise in the system and is proportional to σ^2 . Note that when only one cue is nonneutral, $f_i \propto e^{I_i/\sigma_{\text{eff}}^2}$ ($i = \lambda, v$), and when both cues are nonneutral, $f_{\lambda v} \propto e^{(I_\lambda + I_v)/\sigma_{\text{eff}}^2}$. Therefore, the fractions are related through $f_{\lambda v} \propto f_\lambda \times f_v$, and after normalization, they follow the multiplicative rule (Eq. 1). Fig. 4F shows that Eq. 3 is indeed satisfied by the diffusion model, because the fraction of dominance of percept A obtained from numerical simulations as a function of the total input current (Fig. 4F, blue line) is a sigmoid function (Fig. 4F, red line). This

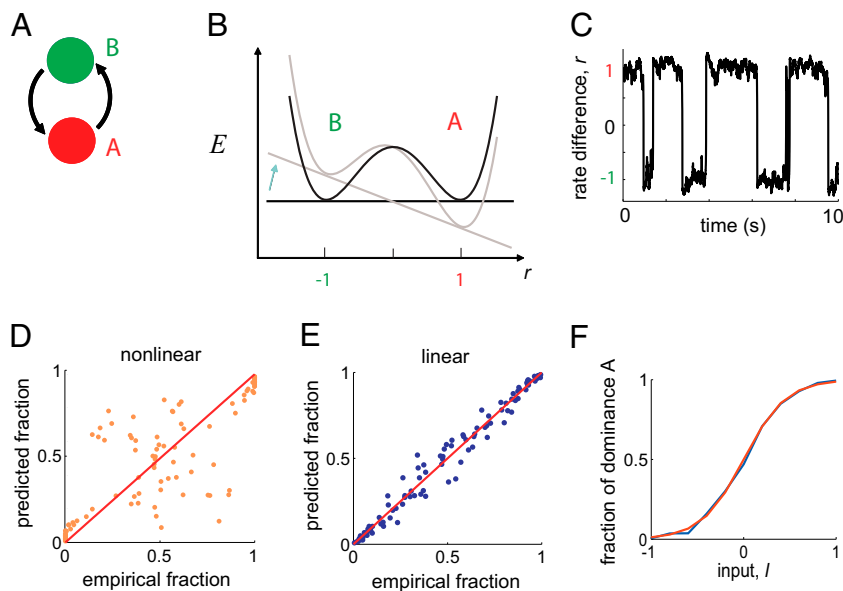


Fig. 4. Simplified network model for Bayesian sampling. (A) Schematic of the neural network. (B) Energy as a function of the difference between the firing rates of the two populations ($r = r_A - r_B$). When the state of the system lies close to the right or left minimum (r is close to 1 or -1), percept A or B dominates, respectively. Alternations in dominance happen because noise can kick the system from one minimum to the other minimum. When the two cues are neutral (black line), the two percepts dominate for equal amounts of time (i.e., $f = 0.5$). When the cues favor percept A, the energy landscape is tilted to the right (gray line), and $f > 0.5$. (C) Population rate difference r as a function of time. Stochastic switches occur between the two states of the system. (D and E) Fractions of dominance predicted by the multiplicative rule vs. observed fractions of dominance generated by the model (D, orange dots and E, blue dots) with nonlinear (D) and linear (E) inputs ($\varepsilon = 5$ and $\varepsilon = 0$, respectively). The model's performance lies close to the unit slope line (red) only when the inputs are combined linearly (E). (F) Fraction of dominance of state A ($r > 0$) as a function of the total input. The curve (blue) is well-fitted by a sigmoid function (red).

analytical approach can also be used to reveal why the system with a nonlinear function does not follow the multiplicative rule. Because in this case, $f_{\lambda\nu} \propto e^{g(I_\lambda, I_\nu)/\sigma_{off}^2}$, the product of the fractions when only one cue is nonneutral is not equal to the fraction when the two cues are nonneutral.

Sampling with Realistic Neural Circuits. The main features of the energy model can be implemented in a neural network with attractor dynamics. We consider a recurrent neural network with two competing populations (Fig. 5A) encoding the two percepts A and B, whose states are described by their population averaged firing rates r_A and r_B , as suggested by neural data (29). An additional relay neuronal population fires in response to the cues and provides inputs to the competing populations A and B with positive (direct connections) and negative (through an inhibitory population) signs, respectively. The firing of the relay population is a function of the sum of the cue strengths, $I_\lambda + I_\nu$. We consider linear and nonlinear activation functions (SI Methods) close to

those functions found in primary visual cortex (25, 26). We also added a slow adaptation process (30–33).

The network stochastically alternates between percepts with gamma-like distributions of dominance durations, which captures several aspects of the experimental distributions (Fig. 5B) (14, 17, 34–36). The distributions generated by the network are not significantly different from those distributions obtained from pooling data across subjects (Fig. 5B) (two-sample Kolmogorov–Smirnov test, $p > 0.05$). The distributions from human data have a coefficient of variation (CV; ratio between SD and mean) close to 0.6, regardless of the fraction of dominance (Fig. 5C, blue dots) (slope not significantly different from zero, $p = 0.3$). Although the model shows a significant linear dependence on the fraction ($p < 0.05$), the dependence is weak, and the CV is consistently close to the experimental value (Fig. 5C, red dots). Importantly, the network predicts that the mean dominance durations of a percept should depend primarily on its fraction of dominance. The experimental data not only show this important qualitative feature but also follow quantitatively the idiosyncratic

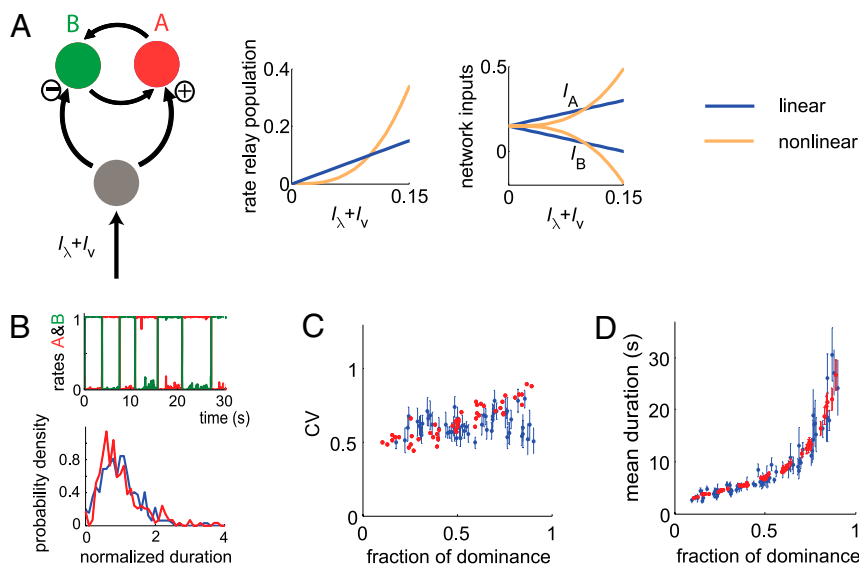


Fig. 5. Sampling and multiplicative rule in attractor neural networks. (A) Architecture of the network, linear, and nonlinear activation functions of the relay population and resulting inputs to the network. (B) Population firing rates as a function of time. (Upper) Red, population A; green, population B. (Lower) Distributions of dominance durations from the neural network model when the cues are neutral (red) and from the pooled data across subjects (blue) for the wavelength speed experiment in the neutral condition ($n = 320$). Time has been normalized so that the mean of the distributions is one. Because the distribution from the model corresponds to the case in which the cues are neutral (zero biasing currents), it is the same regardless of whether the activation function of the relay unit is linear or nonlinear. (C) CV of the dominance duration distribution of a percept as a function of its fraction of dominance for the data averaged across subjects (blue) and model (red). (D) Mean dominance duration of a percept as a function of its fraction of dominance for the experimental data averaged across subjects (blue) and for the model (red). Model error bars correspond to SEM across durations.

mean duration vs. fraction dependence obtained from the model (Fig. 5D). These results hold independently of whether the activation function of the relay population is linear (Fig. 5B–D) or nonlinear (SI Methods and Fig. S2).

The slow dynamics of switches indicate that bistable perception generates temporally correlated samples (successive samples tend to be similar, which is indicated by the fact that percepts tend to linger for hundreds of milliseconds before switching), a property consistent with Langevin Monte Carlo sampling (28).

Therefore, the network generates a stochastic behavior consistent with bistable perception and makes nontrivial predictions about the dynamics of perceptual bistability. However, this behavior does not necessarily mean that the network follows the multiplicative rule. Interestingly, when the activation function in the relay population is nonlinear, the fractions of dominance do not combine multiplicatively (Figs. 2D, 3D, and 6A, orange dots). In contrast, when the activation function is linear-rectified, the network obeys the multiplicative rule (Figs. 2D, 3D, and 6A, blue dots). This result holds because the fraction of dominance time is a sigmoid function of the sum of input currents when the inputs to the network are linear (Fig. 6B, blue lines) but not when the inputs are nonlinear (Fig. 6B, orange lines). We show in SI Methods (Fig. S3) that these results hold even in a more realistic network with integrate and fire neurons.

Probabilistic Population Codes and Bayesian Sampling. We have shown in the previous sections how to build a recurrent network that implements the multiplicative rule, but we have not shown yet that the network samples the posterior distribution over image interpretations specified by the input signals. If the fraction of dominance for a given cue is the result of sampling the posterior distribution over image interpretations $p(s|I_i)$ (here $s = \{A, B\}$ and I_i is the current induced by cue $i = \{\lambda, \nu\}$), then the fraction of dominance and the posterior distribution should be the same function of the input current, I_i . Because the attractor network generates fractions of dominance that are sigmoid functions of the current (Eq. 3), the attractor network is sampling the posterior distribution only if that distribution is also a sigmoid function of the input current, that is (Eq. 4),

$$p(s = A | I_i) = f(s = A | I_i) = \frac{1}{1 + e^{-2I_i/\sigma_{\text{eff}}^2}}. \quad [4]$$

Moreover, through Bayes rule, we know that (Eq. 5)

$$p(s = A | I_i) \propto p(I_i | s = A), \quad [5]$$

where the function $p(I_i | s = A)$ corresponds to the variability in neural responses (in this case, one input current) over multiple presentations of the same stimulus s . Therefore, the key question is whether neural variability in vivo has a distribution consistent with Eqs. 4 and 5. If this is not the case, attractor dynamics would not be sampling from the posterior distributions of s .

Experimentally, neural variability is typically assessed by measuring the variability in spike counts for a fixed s as opposed to the variability in input currents. Mapping input current onto spike counts is easy if we assume, as we did earlier, that the input current is proportional to the difference in spike counts vectors, $\mathbf{r}_A - \mathbf{r}_B$, from two presynaptic populations (e.g., V1 neurons with different depth and speed preferences) (37), one that prefers stimulus $s = A$ and the other that prefers stimulus $s = B$. One can then show (SI Methods) that Eqs. 4 and 5 are only satisfied when the distribution over either \mathbf{r}_A or \mathbf{r}_B given s takes the form $p(\mathbf{r} | s) \propto \phi(\mathbf{r}) \exp(\mathbf{h}(s) \cdot \mathbf{r})$, where $\mathbf{h}(s)$ is a kernel related to the tuning curves and covariance matrix of the neural responses. Remarkably, this family of distributions, known as the exponential family with linear sufficient statistics, provides a very close approximation to the variability observed in vivo (22, 38).

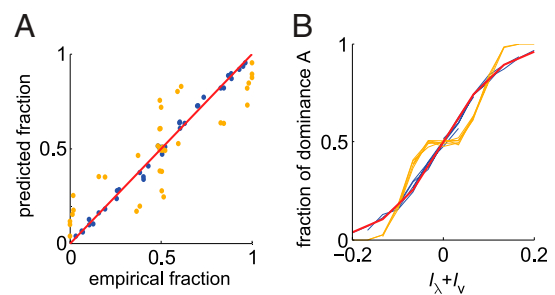


Fig. 6. (A) Predicted fractions from the multiplicative rule vs. observed fractions of dominance generated by the neural network with nonlinear (orange) and linear (blue) inputs (SI Methods). As observed with the energy model (Fig. 4E), the network follows the multiplicative rule only when the relay population has a linear activation function. (B) Fraction of dominance of state A as a function of the total input when the relay population is nonlinear (orange) and linear (blue). The latter are well-fitted by a sigmoid function (red), which was the case with the energy model (Fig. 4F).

This family of distributions corresponds also to a form of neural code known as probabilistic population codes (22). In other words, our results show that attractor dynamics can be used to sample from a posterior distribution encoded by a probabilistic population code using the exponential family with linear sufficient statistics.

Discussion

We have reported that the fraction of dominance in bistable perception behaves as a probability. This result supports the notion that the visual system samples the posterior distribution over image interpretations. In addition, we showed that attractor networks can implement Bayesian sampling only when the variability of neuronal activity follows the exponential family with linear sufficient statistics, as observed experimentally.

This last result is important, but using the exponential family has another advantage. Several works have reported that humans perform near-optimal cue integration in a variety of settings (1–8). It is, therefore, essential that the combination of inputs that leads to the multiplicative rule in an attractor network also results in optimal cue integration. We saw that inputs need to be added to observe the multiplicative rule in an attractor network. Adding two inputs does not necessarily result in optimal cue integration, but again, when the variability of cortical activity follows the exponential family with linear sufficient statistics, it is the optimal combination rule for cue integration (22). Therefore, the fact that the neural variability follows the exponential family allows both Bayesian sampling and optimal integration of evidence with attractor networks.

Our study is not the first study to investigate cue combination and perceptual bistability, but previous works did not test whether bistable perception is akin to what we defined as Bayesian sampling (19, 20). The fact that bistable perception alternates between two interpretations is certainly suggestive of a sampling process but not necessarily of Bayesian sampling. For instance, the orange dots in Fig. 6A show an example of a network that stochastically oscillates with gamma-like distributions over percept durations (Fig. 5B), as observed in our experimental data. The kind of analysis that has been used in previous studies to argue that bistable perception is a form of sampling (19, 20) would also conclude that this network is sampling. However, this particular network does not perform Bayesian sampling; it does not follow the multiplicative rule (Fig. 6A). In contrast, our experimental results make it clear that bistable perception follows the multiplicative rule predicted by Bayesian sampling.

Bayesian sampling has several computational advantages. For instance, in the context of reinforcement learning, when the sta-

tistics of the world is fixed, the optimal solution involves picking the action that is the most likely to be rewarded; however, when the statistics of the world change over the time, sampling from the posterior distribution, which is a form of exploratory behavior (21, 39), is more sensible (40). Interestingly, bistable perception implements a form of sampling that could be used to smoothly interpolate between pure exploration (sampling from the posterior) and pure exploitation (choosing the action that is the most likely to be rewarded). Indeed, our results suggest that bistable perception samples from posterior distributions that are raised to a power, p^n , where n can take any value (*SI Methods*). When n is large, the most likely state is sampled on almost every iteration, which corresponds to exploitation, whereas setting n close to zero leads to exploratory behavior.

The fact that low-level vision and perhaps low-level perception might involve sampling is particularly interesting in light of several other recent findings suggesting that higher-level cognitive tasks, like causal reasoning (41, 42) and decision-making (43), might also involve some form of sampling. Sampling may turn out to be a general algorithm for probabilistic inference in all domains.

Methods

Experimental Methods. The stimulus consisted of two superimposed square-wave gratings, denoted α and β , moving at an angle of 160° between their directions of motion behind a circular aperture (21) (Fig. 1A) with the parameters specified in *SI Methods*. The gratings consisted of gray bars of

equal luminance presented on a white background. Where the gray bars intersected, the luminance was set to that of the bars (as if one of the bars was occluding the other bar). Observers were asked to continually report their percept by holding down one of two designated keys [i.e., motion direction (right or left) of the grating that they perceived as being behind the other grating] and not to press any key if they were not certain. We measured, in each trial, the accumulated time that either percept (i.e., depth ordering) was dominant and computed the fraction of time that percept $s = \{A, B\}$ dominated as $f(s) = (\text{the cumulative time percept } s \text{ was reported as dominant}) / (\text{the total time that either of the percepts was reported as dominant})$. Therefore, this fraction corresponds to the proportion of time that percept s dominated. Percept A denotes the percept in which grating α is behind grating β (and conversely, percept B). Fractions of dominance shown in the figures correspond to averaged values of the fractions across trials and observers, and error bars correspond to SEM across the population.

Mathematical Methods. The derivations of the multiplicative rule and stronger cue take all rule and the descriptions of the energy, rate-based, and spiking models are presented in *SI Methods*.

ACKNOWLEDGMENTS. We thank Jan Drugowitsch and Robbie Jacobs for their suggestions and comments. We are also very grateful to Thomas Thomas and Bo Hu for their assistance during the experimental setup and Vick Rao for his help in using the cluster. D.C.K. is supported by National Institutes of Health Grant EY017939. A.P. is supported by National Science Foundation Grant BCS0446730 and the Multidisciplinary University Research Initiative (MURI) Grant N00014-07-1-0937. This work was also partially supported by National Eye Institute Award P30 EY001319.

- Ernst MO, Banks MS (2002) Humans integrate visual and haptic information in a statistically optimal fashion. *Nature* 415:429–433.
- Jacobs RA (1999) Optimal integration of texture and motion cues to depth. *Vision Res* 39:3621–3629.
- Landy MS, Maloney LT, Johnstone EB, Young M (1995) Measurement and modeling of depth cue combination: In defense of weak fusion. *Vision Res* 35:389–412.
- van Beers RJ, Sittig AC, Gon JJ (1999) Integration of proprioceptive and visual position-information: An experimentally supported model. *J Neurophysiol* 81:1355–1364.
- Körding KP, Wolpert DM (2006) Bayesian decision theory in sensorimotor control. *Trends Cogn Sci* 10:319–326.
- Hillis JM, Watt SJ, Landy MS, Banks MS (2004) Slant from texture and disparity cues: Optimal cue combination. *J Vis* 4:967–992.
- Knill DC (2007) Robust cue integration: A Bayesian model and evidence from cue-conflict studies with stereoscopic and figure cues to slant. *J Vis* 7:5.1–5.24.
- Knill DC (2003) Mixture models and the probabilistic structure of depth cues. *Vision Res* 43:831–854.
- Tjan BS, Braje WL, Legge GE, Kersten D (1995) Human efficiency for recognizing 3-D objects in luminance noise. *Vision Res* 35:3053–3069.
- Gold JM, Tadin D, Cook SC, Blake R (2008) The efficiency of biological motion perception. *Percept Psychophys* 70:88–95.
- Kersten D, Mamassian P, Yuille A (2004) Object perception as Bayesian inference. *Annu Rev Psychol* 55:271–304.
- Hinton GE (2007) Learning multiple layers of representation. *Trends Cogn Sci* 11:428–434.
- Fiser J, Berkes P, Orbán G, Lengyel M (2010) Statistically optimal perception and learning: From behavior to neural representations. *Trends Cogn Sci* 14:119–130.
- Blake R (2001) A primer on binocular rivalry. *Brain and Mind* 2:5–38.
- Blake R, Logothetis NK (2002) Visual competition. *Nat Rev Neurosci* 3:13–21.
- Dayan P (1998) A hierarchical model of binocular rivalry. *Neural Comput* 10:1119–1135.
- Necker LA (1832) Observations on some remarkable phenomenon which occurs on viewing a figure of a crystal of geometrical solid. *Lond Edinburgh Phil Mag J Sci* 3:329–337.
- Rubin E (1958) Figure and ground. *Readings in Perception*, eds Beardslee DC, Wertheimer M (Van Nostrand Reinhold, New York), pp 194–203.
- Sundareswara R, Schrater PR (2008) Perceptual multistability predicted by search model for Bayesian decisions. *J Vis* 8:12.1–12.19.
- Hoyer PO, Hyvarinen A (2003) Interpreting neural response variability as Monte Carlo sampling of the posterior. In Becker S, et al., editors. *Advances in Neural Information Processing Systems* 15. MIT Press; 2003. pp. 277–284.
- Moreno-Bote R, Shpiro A, Rinzel J, Rubin N (2008) Bi-stable depth ordering of superimposed moving gratings. *J Vis* 8:20.1–20.13.
- Ma WJ, Beck JM, Latham PE, Pouget A (2006) Bayesian inference with probabilistic population codes. *Nat Neurosci* 9:1432–1438.
- Mamassian P, Landy MS (1998) Observer biases in the 3D interpretation of line drawings. *Vision Res* 38:2817–2832.
- Mamassian P, Landy MS (2001) Interaction of visual prior constraints. *Vision Res* 41:2653–2668.
- Priebe NJ, Mechler F, Carandini M, Ferster D (2004) The contribution of spike threshold to the dichotomy of cortical simple and complex cells. *Nat Neurosci* 7:1113–1122.
- Carandini M, Ferster D (2000) Membrane potential and firing rate in cat primary visual cortex. *J Neurosci* 20:470–484.
- Moreno-Bote R, Rinzel J, Rubin N (2007) Noise-induced alternations in an attractor network model of perceptual bistability. *J Neurophysiol* 98:1125–1139.
- Bishop CM (2006) *Pattern Recognition and Machine Learning* (Springer, Berlin).
- Sheinberg DL, Logothetis NK (1997) The role of temporal cortical areas in perceptual organization. *Proc Natl Acad Sci USA* 94:3408–3413.
- Shpiro A, Moreno-Bote R, Rubin N, Rinzel J (2009) Balance between noise and adaptation in competition models of perceptual bistability. *J Comput Neurosci* 27:37–54.
- Laing CR, Chow CC (2002) A spiking neuron model for binocular rivalry. *J Comput Neurosci* 12:39–53.
- Markram H, Tsodyks M (1996) Redistribution of synaptic efficacy between neocortical pyramidal neurons. *Nature* 382:807–810.
- Abbott LF, Varela JA, Sen K, Nelson SB (1997) Synaptic depression and cortical gain control. *Science* 275:220–224.
- Leopold DA, Logothetis NK (1996) Activity changes in early visual cortex reflect monkeys' percepts during binocular rivalry. *Nature* 379:549–553.
- Levelt WJM (1968) *On Binocular Rivalry* (Mouton, The Hague).
- Hupé JM, Rubin N (2003) The dynamics of bi-stable alternation in ambiguous motion displays: A fresh look at plaids. *Vision Res* 43:531–548.
- Cumming BG, DeAngelis GC (2001) The physiology of stereopsis. *Annu Rev Neurosci* 24:203–238.
- Graf AB, Kohn A, Jazayeri M, Movshon JA (2011) Decoding the activity of neuronal populations in macaque primary visual cortex. *Nat Neurosci* 14:239–245.
- Moreno-Bote R, Shpiro A, Rinzel J, Rubin N (2010) Alternation rate in perceptual bistability is maximal at and symmetric around equi-dominance. *J Vis* 10(11):1,1–18.
- Sutton RS, Barto AG (1998) Reinforcement learning: An introduction. *Adaptive Computation and Machine Learning* (MIT Press, Cambridge, MA).
- Griffiths TL, Tenenbaum JB (2005) Structure and strength in causal induction. *Cognit Psychol* 51:334–384.
- Tenenbaum JB, Griffiths TL, Kemp C (2006) Theory-based Bayesian models of inductive learning and reasoning. *Trends Cogn Sci* 10:309–318.
- Sugrue LP, Corrado GS, Newsome WT (2004) Matching behavior and the representation of value in the parietal cortex. *Science* 304:1782–1787.