

Supplementary information

Fluctuating internal states mediate neural-behavioral covariations in V1

Baowang Li, Jason Samonds, Yuzhi Chen,
Thibaud Taillefumier, Nicholas Priebe, and Eyal Seidemann

1 Decision models

1.1 Additive model

We denote by V_i the membrane voltage of the recorded neuron whose receptive field location determines the “in side” of the field of view. Similarly, we denote by V_o the membrane voltage of a neuron with symmetric receptive-field location in the opposite hemifield, in the “out side” of the field of view. In the additive model, we consider that these voltages are stochastically related via

$$\begin{aligned} V_i &= X_i + A_i + N_{i,e}, \\ V_o &= X_o + A_o + N_{o,e}, \end{aligned}$$

where X_i and X_o denote stimulus strengths, $N_{i,e}$ and $N_{o,e}$ denote (early) sensory noise. We further consider that variable internal states contribute to voltage fluctuations as nonnegative additive random variables A_i and A_o . Assume a distribution p_A such that $A \geq 0$ with $\mathbb{E}[A] = 1$. This internal state can be thought as representing attention or arousal levels, and for conciseness, we will refer to A_i and A_o as attention variables.

Given a criterion c , the decision variable is $D = \mathbb{1}_{\{V_i - V_o \geq c + N_d\}}$, where N_d denotes some hemisphere-specific downstream noise. All noise distributions will be assumed Gaussian centered, which means that we need to consider two standard deviation parameters σ_e and σ_d .

1.2 Multiplicative model

As for the additive model, V_i and V_o refers to the voltage on the “in side” and the “out side”, respectively. In the multiplicative model, we consider that these voltages follows

$$\begin{aligned} V_i &= A_i(X_i + N_{i,e}) + N_{i,l}, \\ V_o &= A_o(X_o + N_{o,e}) + N_{o,l}, \end{aligned}$$

where X_i and X_o denote stimulus strengths, $N_{i,e}$ and $N_{o,e}$ denote early sensory noise, and $N_{i,l}$ and $N_{o,l}$ denote late sensory noise. The hallmark of the multiplicative model is to consider that variable internal states are modeled by gain random variables A_i and A_o , with some distribution p_A such that $A \geq 0$ with $\mathbb{E}[A] = 1$. As for the additive model, this internal state can be thought as representing attention or arousal levels, and for conciseness, we will refer to A_i and A_o as attention variables. One could allow for gain variables to also contribute additively to voltage fluctuations by including constant offset terms within parentheses in the above equations. Such inclusions, however, will not impact our results, and for this reason, we only consider A_i and A_o as pure gain variables in the multiplicative model. Given some criterion c , the decision variable is $D = \mathbb{1}_{\{V_i - V_o \geq c + N_{i,d} + N_{o,d}\}}$, where $N_{i,d}$ and $N_{o,d}$ denote some hemisphere-specific downstream noise. All noise distributions will be assumed independent Gaussian centered, which means that we need to consider three standard deviation parameters σ_e , σ_l , and σ_d .

1.3 Probabilities of decision outcomes

For simplicity, we will always assume no decision bias, i.e. $c = 0$. There is no lack of generality in doing so as considering a decision bias is equivalent to shifting the means μ_{V_i} and μ_{V_o} , which features as free parameters in our derivations. The probability of the decision outcomes D_i and D_o admits an explicit integral expression involving the joint probabilities of the voltages $p(v_i, v_o | (X_i, X_o))$. To see this, first observe that assuming knowledge of (X_i, X_o) and denoting the decision variable as $W = V_i - V_o + N_{i,d} - N_{o,d}$, we have

$$p(W, V_i, V_o) = p(W | V_i, V_o)p(V_i, V_o) = g_{\sqrt{2}\sigma_d}(W - (V_i - V_o))p(V_i, V_o),$$

by independence of the hemisphere-specific decision noises. Then, we have

$$\begin{aligned} p[D_i | (X_i, X_o)] &= \int dv_i \int dv_o \int dw \mathbb{1}_{\{w \geq 0\}} p(w, v_i, v_o | (X_i, X_o)), \\ &= \int dv_i \int dv_o p(v_i, v_o | (X_i, X_o)) \int dw \mathbb{1}_{\{w \geq 0\}} g_{\sqrt{2}\sigma_d}(w - (v_i - v_o)), \\ &= \int dv_i \int dv_o p(v_i, v_o | (X_i, X_o)) (1 + \operatorname{erf}[(v_i - v_o)/\sigma_d]) / 2, \end{aligned}$$

so that we also have

$$p[D_o | (X_i, X_o)] = 1 - p[D_i | (X_i, X_o)] = \int dv_i \int dv_o p(v_i, v_o | (X_i, X_o)) (1 - \operatorname{erf}[(v_i - v_o)/\sigma_d]) / 2.$$

where $\operatorname{erf}[\cdot]$ denotes the error function.

1.4 Symmetric probability model for variable internal state

For simplicity, we consider a model with variable internal states for which each hemisphere alternates between a low-attention state $A_{ll} > 0$ and a high-attention state $A_{hh} > 0$ with $A_{ll} < A_{hh}$. To probe the role of correlations across hemispheres while maintaining symmetry, we consider that the attention states of both hemispheres can be recapitulated by the following two-dimensional probabilistic model:

$$\begin{aligned} \mathbb{P}[A_{ll}] &= \mathbb{P}[A_i = A_{ll}, A_o = A_{ll}] = p, & \mathbb{P}[A_{hl}] &= \mathbb{P}[A_i = A_{hh}, A_o = A_{ll}] = q, \\ \mathbb{P}[A_{lh}] &= \mathbb{P}[A_i = A_{ll}, A_o = A_{hh}] = q, & \mathbb{P}[A_{hh}] &= \mathbb{P}[A_i = A_{hh}, A_o = A_{hh}] = 1 - p - 2q. \end{aligned}$$

For $q = 0$, there is perfect correlation across hemispheres, i.e., (A_i, A_o) alternates between (A_{ll}, A_{ll}) and (A_{hh}, A_{hh}) with probability p , whereas for $q = 1/2$, which implies $p = 0$, there is perfect anticorrelation across hemispheres, i.e., $A_i = A_{hh}$ whenever $A_o = A_{ll}$ and vice versa. More generally, the correlation across hemispheres satisfies

$$\rho = \frac{p - (p + q)^2}{(p + q)(1 - p - q)}, \quad (1)$$

so that zero correlation is achieved for $q = \sqrt{p}(1 - \sqrt{p}) \leq 1/4$, with independence for $q = p = 1/4$.

1.5 Choice discriminability via δ -metric

Here, we assume that the gain variables are fixed to deterministic values for each hemisphere: $A_i = a_i$ and $A_o = a_o$. Under these assumptions, the models under consideration are linear, allowing one to carry out explicit calculations for Gaussian noise. Specifically, one can find a closed-form expression for a simple

choice-discriminability metric associated to the alternative forced choice D . Denoting the “choice in” by $D_i = D$ and the “choice out” by $D_o = 1 - D$ for conciseness, this metric is defined as a difference between conditional expectations

$$\delta(T) = \mathbb{E}[V_i | D_i, T] - \mathbb{E}[V_i | D_o, T] = \frac{\mathbb{E}[V_i D_i | T]}{\mathbb{P}[D_i | T]} - \frac{\mathbb{E}[V_i D_o | T]}{\mathbb{P}[D_o | T]},$$

where the conditioning is with respect to the decision outcome D . This metric is closely related to the classical δ -metric, which is derived from δ by normalizing with respect to the standard deviations of the involved distributions. We consider the δ -metric for three stimulus conditions. In all these conditions, we assume that the presence of target-in the receptive field of a neuron leads to an elevation of the mean voltage level for this neuron, i.e., the presence or absence of a stimulus corresponds to values x_1 and x_0 , respectively, such that $x_1 > x_0$. Then, the three considered stimulus conditions are:

“target-in”: $T_i = \{X_i = x_1, X_o = x_0\}$.

“target-out”: $T_o = \{X_i = x_0, X_o = x_1\}$.

“no-stimulus”: $T_n = \{X_i = x_0, X_o = x_0\}$.

2 Choice discriminability with fixed attention state

Independent of the stimulus condition, given fixed attention states $A_i = a_i$ and $A_o = a_o$, V_i and V_o follow independent Gaussian distributions so that

$$p(V_i, V_o | (X_i, X_o)) = g_{\sigma_{V_i}}(v_i - \mu_{V_i}) g_{\sigma_{V_o}}(v_o - \mu_{V_o}),$$

where $g_{\sigma_{V_i}}$ and $g_{\sigma_{V_o}}$ denotes the density of centered Gaussian variables with variance $\sigma_{V_i}^2$ and $\sigma_{V_o}^2$, respectively. The means and variances are given by:

additive model : $\mu_{V_i} = X_i + a_i, \mu_{V_o} = X_o + a_o, \sigma_{V_i}^2 = \sigma_e^2, \sigma_{V_o}^2 = \sigma_e^2$.

multiplicative model: $\mu_{V_i} = a_i X_i, \mu_{V_o} = a_o X_o, \sigma_{V_i}^2 = a_i^2 \sigma_e^2 + \sigma_1^2, \sigma_{V_o}^2 = a_o^2 \sigma_e^2 + \sigma_1^2$.

For a fixed attention state, independent of the nature of the model, the decision process based on a competition between the voltage signals of both hemispheres can be represented schematically. We utilize such a representation in Fig 1 in relation to the definition of the δ -metric as the difference between two choice-conditioned voltage means for V_i . In the Gaussian context, one can turn this schematic view into analytical expressions allowing one to explore the parametric dependence of the δ -metric. For instance, one can express the probability of the decision outcomes D_i and D_o explicitly. The derivation relies on performing the change of variable $(x, y) = (v_i - v_o, v_i + v_o)$ to evaluate two iterated integrals. For instance,

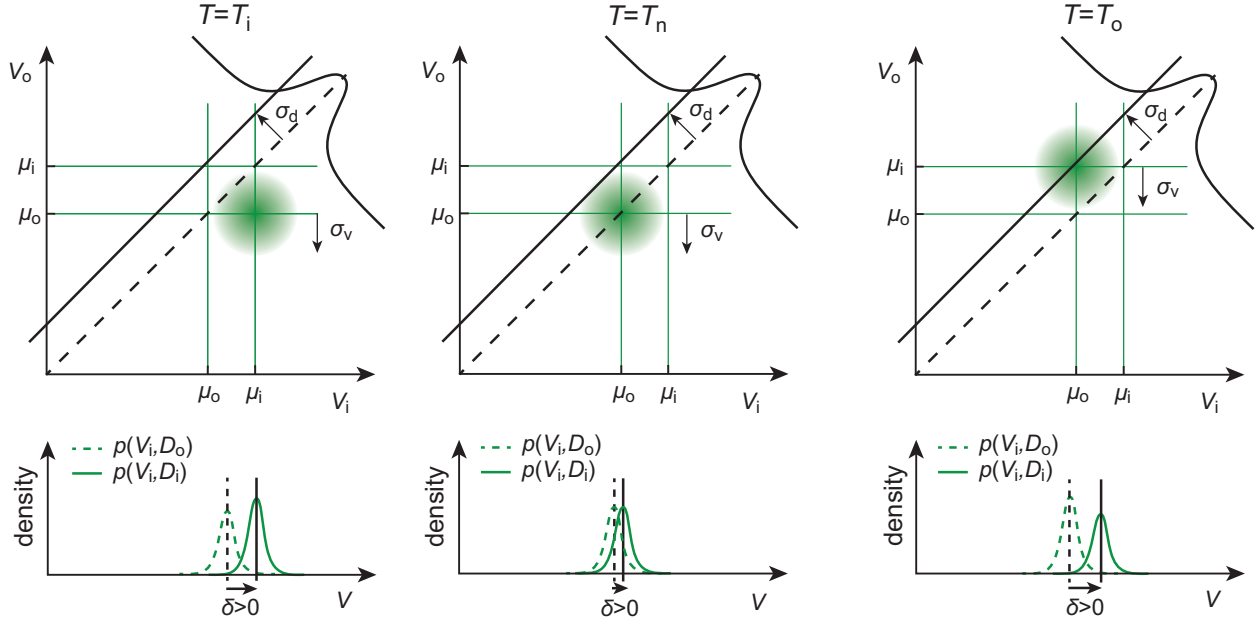


Figure 1: **Schematic for decision with fixed gain.** We consider three stimulus conditions: “target-in”, i.e., $T = T_i$, “no-stimulus”, i.e., $T = T_n$, and “target-out”, i.e., $T = T_o$. For fixed gain, the model is purely Gaussian for all stimulus conditions, which impact the distribution of voltages (V_i, V_o) via variables means. The assumption is that the presence of a stimulus is assumed to cause an increase in mean voltages as depicted in the two-dimensional phase space (top panels). Within this picture, the decision outcome is determined by the position of the voltage point (V_i, V_o) with respect to the decision boundary, a slope-one line with Gaussian-distributed intercept with variance σ_d^2 (solid black line). Recorded voltage data correspond to the projection of the phase space on the V_i -space conditioning to the decision outcome (bottom panels): for $D = D_i$ and $D = D_o$, this corresponds to marginalizing the distribution of voltages (V_i, V_o) restricted to the portion of the phase plane below or above the fluctuating decision boundary (solid black line in top panels), respectively. For large downstream decision variance, the joint probability densities $p(V_i, D_i | T)$ and $p(V_i, D_o | T)$ should integrate to about the same choice probabilities, leading to nonnegative but low δ -metric values.

the probability of “choice in” is given as

$$\begin{aligned}
\mathbb{P}[D_i | (X_i, X_o)] &= \frac{1}{4\pi\sigma_{V_i}\sigma_{V_o}} \iint \exp\left(-\frac{(v_i - \mu_{V_i})^2}{2\sigma_{V_i}^2} - \frac{(v_o - \mu_{V_o})^2}{2\sigma_{V_o}^2}\right) \left(1 + \operatorname{erf}\left[\frac{v_i - v_o}{\sigma_d}\right]\right) dv_i dv_o, \\
&= \frac{1}{2\pi\sigma_{V_i}\sigma_{V_o}} \iint \exp\left(-\frac{(x + y - \mu_{V_i})^2}{2\sigma_{V_i}^2} - \frac{(y - x - \mu_{V_o})^2}{2\sigma_{V_o}^2}\right) \left(1 + \operatorname{erf}\left[\frac{x}{\sigma_d}\right]\right) dx dy, \\
&= \frac{1}{\sqrt{2\pi}\sqrt{\sigma_{V_i}^2 + \sigma_{V_o}^2}} \int \exp\left(-\frac{(2x + \mu_{V_o} - \mu_{V_i})^2}{2(\sigma_{V_i}^2 + \sigma_{V_o}^2)}\right) \left(1 + \operatorname{erf}\left[\frac{x}{\sigma_d}\right]\right) dx, \\
&= \frac{1}{2} \left(1 + \operatorname{erf}\left[\frac{\mu_{V_i} - \mu_{V_o}}{\sqrt{2}(\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2)}\right]\right),
\end{aligned}$$

where the last equality follows from the identity

$$\int \sqrt{\frac{a}{\pi}} e^{-(ax+b)^2} \operatorname{erf}[x] dx = -\operatorname{erf}\left[\frac{b}{\sqrt{a^2 + 1}}\right].$$

The complementary probability of “choice out” directly follows as

$$\mathbb{P}[D_o | (X_i, X_o)] = 1 - \mathbb{P}[D_i | (X_i, X_o)] = \frac{1}{2} \left(1 - \operatorname{erf}\left[\frac{\mu_{V_i} - \mu_{V_o}}{\sqrt{2}(\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2)}\right]\right).$$

Then, one can use similar calculations to derive the closed-form expression of the expectations $\mathbb{E}[V_i D_i | T]$ and $\mathbb{E}[V_i D_o | T]$. For instance, we find that

$$\begin{aligned}
\mathbb{E}[V_i D_i | (X_i, X_o)] &= \frac{1}{4\pi\sigma_{V_i}\sigma_{V_o}} \iint v_i \exp\left(-\frac{(v_i - \mu_{V_i})^2}{2\sigma_{V_i}^2} - \frac{(v_o - \mu_{V_o})^2}{2\sigma_{V_o}^2}\right) \left(1 + \operatorname{erf}\left[\frac{v_i - v_o}{\sigma_d}\right]\right) dv_i dv_o, \\
&= \frac{1}{4\pi\sigma_{V_i}\sigma_{V_o}} \iint (x + y) \exp\left(-\frac{(x + y - \mu_{V_i})^2}{2\sigma_{V_i}^2} - \frac{(y - x - \mu_{V_o})^2}{2\sigma_{V_o}^2}\right) \left(1 + \operatorname{erf}\left[\frac{x}{\sigma_d}\right]\right) dv_i dv_o, \\
&= \iint \frac{\mu_{V_o}\sigma_{V_i}^2 + \mu_{V_i}\sigma_{V_o}^2 + 2\sigma_{V_i}^2 x}{\sqrt{2\pi}(\sigma_{V_i}^2 + \sigma_{V_o}^2)^{3/2}} \exp\left(-\frac{(2x + \mu_{V_o} - \mu_{V_i})^2}{2(\sigma_{V_i}^2 + \sigma_{V_o}^2)}\right) \left(1 + \operatorname{erf}\left[\frac{x}{\sigma_d}\right]\right) dv_i dv_o, \\
&= \frac{\mu_{V_i}}{2} \left(1 + \operatorname{erf}\left[\frac{\mu_{V_i} - \mu_{V_o}}{\sqrt{2}(\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2)}\right]\right) + \\
&\quad \frac{\sigma_{V_i}^2}{\sqrt{2\pi}\sqrt{\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2}} \exp\left(-\frac{(\mu_{V_i} - \mu_{V_o})^2}{2(\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2)}\right).
\end{aligned}$$

Normalizing by the choice probability $\mathbb{P}[D_i | (X_i, X_o)]$ yields the conditional expectation

$$\mathbb{E}[V_i | D_i, (X_i, X_o)] = \mu_{V_i} + \frac{\sqrt{2}\sigma_{V_i}^2}{\sqrt{\pi}\sqrt{\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2}} \left(\frac{\exp\left(-\frac{(\mu_{V_i} - \mu_{V_o})^2}{2(\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2)}\right)}{1 + \operatorname{erf}\left[\frac{\mu_{V_i} - \mu_{V_o}}{\sqrt{2}(\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2)}\right]} \right),$$

whereas similar calculations for the “choice out” yield

$$\mathbb{E}[V_i | D_o, (X_i, X_o)] = \mu_{V_i} - \frac{\sqrt{2}\sigma_{V_i}^2}{\sqrt{\pi}\sqrt{\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2}} \left(\frac{\exp\left(-\frac{(\mu_{V_i} - \mu_{V_o})^2}{2(\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2)}\right)}{1 - \operatorname{erf}\left[\frac{\mu_{V_i} - \mu_{V_o}}{\sqrt{2}\sqrt{\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2}}\right]} \right).$$

As expected, conditioning on the decision outcome D introduces biases in the evaluation of the expected values of V_i and V_o , which otherwise evaluate to μ_{V_i} and μ_{V_o} , respectively, in the absence of conditioning. One can check that these biases vanish in the limit of large downstream decision noise $\sigma_d \rightarrow \infty$, as the decision no longer depends on the voltages V_i and V_o . Moreover, these biases simplify when $b = \mu_{V_i} - \mu_{V_o}$ is large compared with the total noise level. Indeed, when $(\mu_{V_i} - \mu_{V_o})^2 \gg \sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2$, we have

$$\begin{aligned} \mathbb{E}[V_i | D_i, (X_i, X_o)] &\underset{b \rightarrow \infty}{\sim} \mu_{V_i} \quad \text{and} \quad \mathbb{E}[V_i | D_o, (X_i, X_o)] \underset{b \rightarrow \infty}{\sim} \mu_{V_i} \left(1 - \frac{\sigma_{V_i}^2}{\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2}\right), \\ \mathbb{E}[V_i | D_o, (X_i, X_o)] &\underset{b \rightarrow -\infty}{\sim} \mu_{V_i} \quad \text{and} \quad \mathbb{E}[V_i | D_i, (X_i, X_o)] \underset{b \rightarrow -\infty}{\sim} \mu_{V_i} \left(1 - \frac{\sigma_{V_i}^2}{\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2}\right). \end{aligned}$$

Thus conditioning gets more impactful when most of the variability originates from σ_{V_i} as opposed to σ_{V_o} and σ_{V_d} .

Finally, the corresponding δ -metric admits the following closed-form expression

$$\delta(T) = \frac{2\sqrt{2}\sigma_{V_i}^2}{\sqrt{\pi}\sqrt{\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2}} \left(\frac{\exp\left(-\frac{(\mu_{V_i} - \mu_{V_o})^2}{2(\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2)}\right)}{1 - \left(\operatorname{erf}\left[\frac{\mu_{V_i} - \mu_{V_o}}{\sqrt{2}\sqrt{\sigma_{V_i}^2 + \sigma_{V_o}^2 + 2\sigma_d^2}}\right]\right)^2} \right),$$

where the involved means μ_{V_i} and μ_{V_o} and variances $\sigma_{V_i}^2$ and $\sigma_{V_o}^2$ depend on the stimulus condition and on the decision model. Interestingly, the above expression is symmetric with respect to μ_{V_i} and μ_{V_o} . Thus, if the noise levels are identical across hemisphere, i.e., if the variances are independent of the stimulus condition such that $\sigma_{V_i}^2 = \sigma_{V_o}^2 = \sigma_V^2$, we expect the δ -metric to be identical for “choice in” and “choice out”. By contrast, the δ -metric attains its minimum in the no-stimulus condition for which we have $\mu_{V_i} = \mu_{V_o}$, which leads to

$$\delta(T_n) = \frac{2\sigma_V}{\sqrt{\pi}\sqrt{1 + (\sigma_d/\sigma_V)^2}} < \delta(T_i) = \delta(T_o).$$

3 Choice discriminability with fluctuating attention state

3.1 Heuristic picture

We explain the heuristics of our approach by first considering the additive model with perfect correlation in Fig. 3. For perfect correlation, $q = 0$ so that either both hemispheres are in the low attention state, i.e., $(A_i, A_o) = (A_{ll}, A_{ll})$, a state we denote by A_{ll} as a shorthand, or both hemispheres are in the high attention state, i.e., $(A_i, A_o) = (A_{hh}, A_{hh})$, a state we denote by A_{hh} as a shorthand. Let us further assume that $p = 1/2$ for simplicity so that the low-attention state and the high-attention state are equiprobable. In the

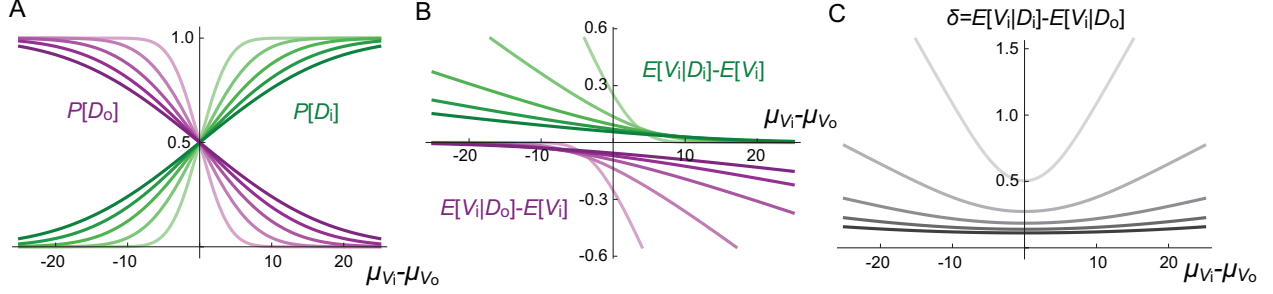


Figure 2: **Decision statistics for the deterministic Gaussian gain model.** **A.** Probabilities of decision as functions of the unconditional mean difference $\mu_{V_i} - \mu_{V_o}$. **B.** Conditioning bias as a function of the unconditional mean difference $\mu_{V_i} - \mu_{V_o}$. **C.** δ -metric as a function of the unconditional mean difference $\mu_{V_i} - \mu_{V_o}$. For large downstream decision variance σ_d^2 , the δ -metric becomes essentially independent of the stimulus-driven mean activity as specified by $\mu_{V_i} - \mu_{V_o}$. Parameters: $\mu_{V_i} = 0$, $\sigma_V^2 = 1$, $\sigma_d^2 = 1, 2, \dots, 5$, increasing variance corresponds to increasing downstream decision variance.

additive model, the effect of changing attention levels is to shift the voltage distribution of (V_i, V_o) along the slope-one decision boundary. Then, by symmetry, marginalizing over V_o given the low-attention state or the high-attention state leads to two conditional distributions $p(V_i, D_i | A_{ll}, T_i)$ and $p(V_i, D_o | A_{ll}, T_o)$ that are translated version of one another. This is consistent with the fact that choices are made on the basis of $V_i - V_o$, i.e., along a direction that is orthogonal to the decision boundary. Finally, the δ -metric results from averaging the conditional distributions $p(V_i, D_i | A_{ll}, T_i)$ and $p(V_i, D_o | A_{ll}, T_o)$ over the equiprobable attention states. Such averaging leaves the δ -metric metric unchanged compared to the case of a single, fixed attention state.

One could have spared the arguments given above by simply remarking that the decision variable cancels out additive variables that are perfectly correlated across hemispheres: if $A_i = A_o = A$, we have $D = \mathbb{1}_{W \geq 0}$ with $W = (V_i + A + N_{i,d}) - (V_i + A + N_{o,d}) = (V_i + N_{i,d}) - (V_i + N_{o,d})$. Therefore, modeling gain as a purely additive shared component across hemispheres cannot alter decision making. The main interest for making the arguments developed above is to consider them in relation to the multiplicative model with perfect correlation depicted in Fig. 4. In the low-attention state A_{ll} , the biased sensory input that follows the presentation of a stimulus is not large enough to significantly overcome the downstream noise, leading to approximately equiprobable decision outcomes, similarly as for the additive model in Fig. 3. This corresponds to the fact that the densities $p(V_i, D_o | A_{ll}, T)$ and $p(V_i, D_i | A_{ll}, T)$ in Fig. 4 both approximately integrate to probability one half, independent of the stimulus T . By contrast, in the high-attention state A_{hh} , the amplification of sensory input is sufficient to significantly bias decision toward the correct outcome. This is schematically illustrated by the fact that in the presence of a stimulus, the density $p(V_i, D_i | A_{hh}, T_i)$ and $p(V_i, D_o | A_{hh}, T_o)$, which correspond to hits, integrate to a much larger probability than the density $p(V_i, D_o | A_{ll}, T_i)$ and $p(V_i, D_i | A_{ll}, T_o)$, which corresponds to misses.

Following classical Bayesian analysis, the imbalance in choice probabilities in the the high-attention state A_{hh} compared to the approximately equiprobable probabilities in the low attention state A_{ll} can significantly impact the evaluation of decision metrics such as the δ -metric. For instance, if one assume that low- and high-attention states are equiprobable as in Fig. 4, one can qualitatively explain how the δ -metric can become negative in the “target-out” condition. The key is to recognize that the overall choice probability densities $p(V_i, D_i | T_o)$ and $p(V_i, D_o | T_o)$ are obtained by mixing the gain specific densities with equal contributions. As a result of this mixture, and because the miss probability is low in the high-attention state compared to the low attention state, the vast majority of “choices in” occurs in the low attention state, so

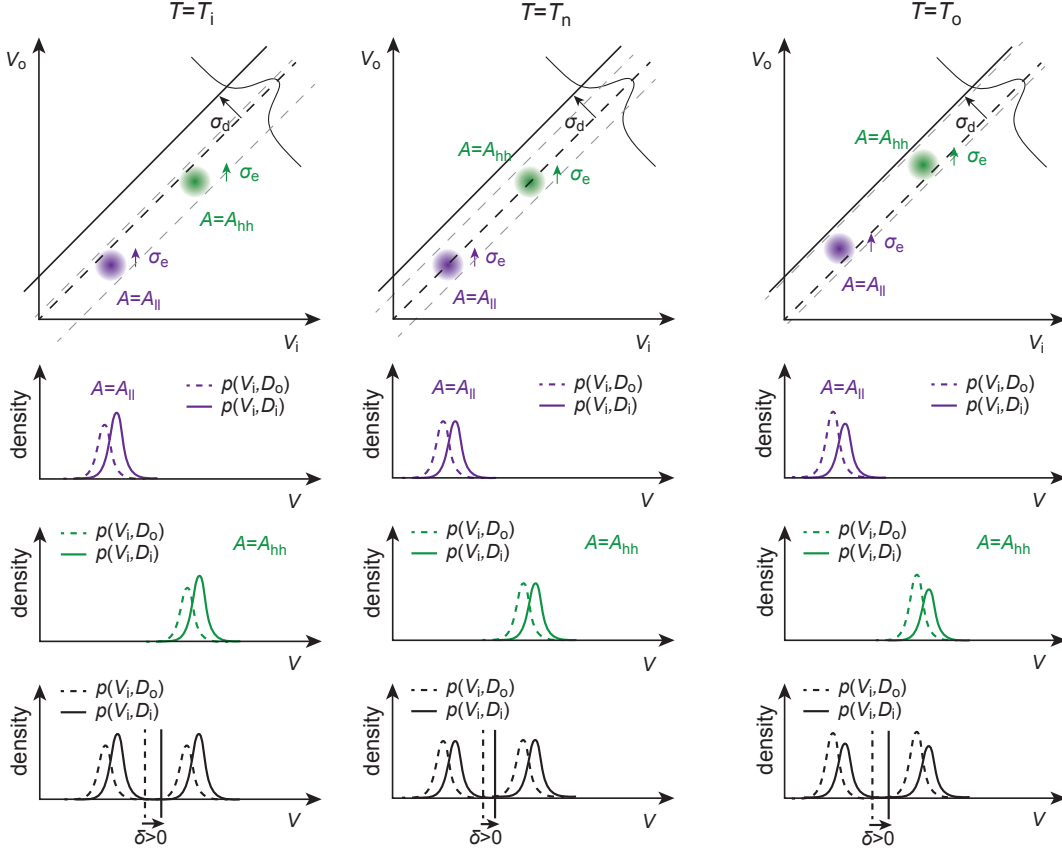


Figure 3: **Schematic for the additive decision model with perfectly correlated fluctuating gain.** We consider three stimulus conditions: “target-in”, i.e., $T = T_i$, “no-stimulus”, i.e., $T = T_n$, and “target-out”, i.e., $T = T_o$; and two equiprobable perfectly correlated attention states: “low gain”, i.e., $A = A_{||}$, and “high gain”, i.e., $A = A_{hh}$. The effect of varying attention states is simply to jointly shift the mean values of voltage values (V_i , V_o) along the decision boundary in the two-dimensional phase space (top panels). When conditioned on an attention state, either $A_{||}$ or A_{hh} , the model is identical to the case of fixed attention state in Fig. 1. Given an attention state, the recorded voltage data can be obtained by the same marginalization operation of the distribution of (V_i , V_o) restricted to each decision region of the two-dimensional phase space as in Fig. 1. Because the δ -metric is invariant with respect to translation along the decision boundary (solid black lines in the top panels), the conditional δ -metric is the same for both attention state in each stimulus condition. The full unconditioned δ -metric is obtained as an average of the two (identical) conditional metric. Thus, considering an additive fluctuating attention states have no impact on the δ -metric compared with the fixed attention state.

that the conditional expectation $\mathbb{E}[V_i | D_i, T_o]$ is only marginally larger than $\mathbb{E}[V_i | D_i, A_{ll}, T_o]$. By contrast, the majority of “choices out” happens in the high-attention state, so that the conditional expectation $\mathbb{E}[V_i | D_o, T_o]$ can considerably exceed $\mathbb{E}[V_i | D_o, A_{ll}, T_o]$, and even approaches $(2\mathbb{E}[V_i | D_o, A_{hh}, T_o] + \mathbb{E}[V_i | D_o, A_{ll}, T_o])/3$ when the high-attention miss probability is close to zero. In the latter extreme scenario, we expect to have

$$\begin{aligned}\delta(T_o) &= \mathbb{E}[V_i | D_i, T_o] - \mathbb{E}[V_i | D_o, T_o] , \\ &\simeq \mathbb{E}[V_i | D_i, A_{ll}, T_o] - (2\mathbb{E}[V_i | D_o, A_{hh}, T_o] + \mathbb{E}[V_i | D_o, A_{ll}, T_o])/3 , \\ &\simeq \frac{1}{3}\delta(T_o, A_{ll}) - \frac{2}{3}(\mathbb{E}[V_i | D_o, A_{hh}, T_o] - \mathbb{E}[V_i | D_i, A_{ll}, T_o]) ,\end{aligned}$$

where we note that $\delta(T_o, A_{ll})$ is typically small in the low-attention state. Thus, $\delta(T_o)$ is negative whenever the high-attention state amplification is strong enough so that the voltage expectation given D_o and A_{hh} is larger than the voltage expectation given D_i and A_{ll} .

3.2 Quantitative argument

Here, we turn our heuristic arguments into quantitative ones for the proposed symmetric model with fluctuating gain level and varying degree of correlation across hemispheres. By symmetry, we denote the miss probability conditioning on the attention state by

$$\begin{aligned}m_{ll} &= \mathbb{P}[D_o | A_{ll}, T_i] = \mathbb{P}[D_i | A_{ll}, T_o] , \\ m_{lh} &= \mathbb{P}[D_o | A_{lh}, T_i] = \mathbb{P}[D_i | A_{hl}, T_o] , \\ m_{hl} &= \mathbb{P}[D_o | A_{hl}, T_i] = \mathbb{P}[D_i | A_{lh}, T_o] , \\ m_{hh} &= \mathbb{P}[D_o | A_{hh}, T_i] = \mathbb{P}[D_i | A_{hh}, T_o] .\end{aligned}$$

All the above miss probabilities can be evaluated explicitly by the same treatment as for the case of a fixed, single attention state. Our goal is to evaluate explicitly the δ -metric in the above fluctuating setting. To do so, we observe that

$$\mathbb{E}[V_i | D, T] = \frac{\mathbb{E}[V_i D | T]}{\mathbb{P}[D | T]} = \sum_{\alpha, \beta \in \{l, h\}} \frac{\mathbb{E}[V_i D, A_{\alpha\beta} | T]}{\mathbb{P}[D | T]} = \frac{1}{\mathbb{P}[D | T]} \sum_{\alpha, \beta \in \{l, h\}} \mathbb{E}[V_i D | A_{\alpha\beta}, T] \mathbb{P}[A_{\alpha\beta}] ,$$

where we assume that the gain fluctuates independently of the stimulus condition: $\mathbb{P}[A | T] = \mathbb{P}[A]$. All the expectations featured above can be evaluated explicitly by the same treatment as for the case of a fixed, single attention state. Thus, it only remains to evaluate the conditional probabilities $\mathbb{P}[D | T]$. In the presence of a stimulus, this involves estimating the marginal miss probability m , satisfying $m = \mathbb{P}[D_o | T_i] = \mathbb{P}[D_i | T_o] = (1 - \mathbb{P}[D_i | T_i]) = (1 - \mathbb{P}[D_o | T_o])$ with

$$m = \sum_{\alpha, \beta \in \{l, h\}} \mathbb{P}[D_o | A_{\alpha\beta}, T_i] \mathbb{P}[A_{\alpha\beta} | T_i] = m_{ll}p + (m_{hl} + m_{lh})q + m_{hh}(1 - p - 2q) ,$$

whereas in the absence of a stimulus, we have $\mathbb{P}[D_i | T_n] = \mathbb{P}[D_o | T_n] = 1/2$ by symmetry. From there, we can express explicitly the δ -metric for all stimulus conditions via the formulas

$$\delta(T) = \sum_{\alpha, \beta \in \{l, h\}} \left(\mathbb{E}[V_i D_i | A_{\alpha\beta}, T] \frac{\mathbb{P}[A_{\alpha\beta}]}{\mathbb{P}[D_i | T]} - \mathbb{E}[V_i D_o | A_{\alpha\beta}, T] \frac{\mathbb{P}[A_{\alpha\beta}]}{\mathbb{P}[D_o | T]} \right) .$$

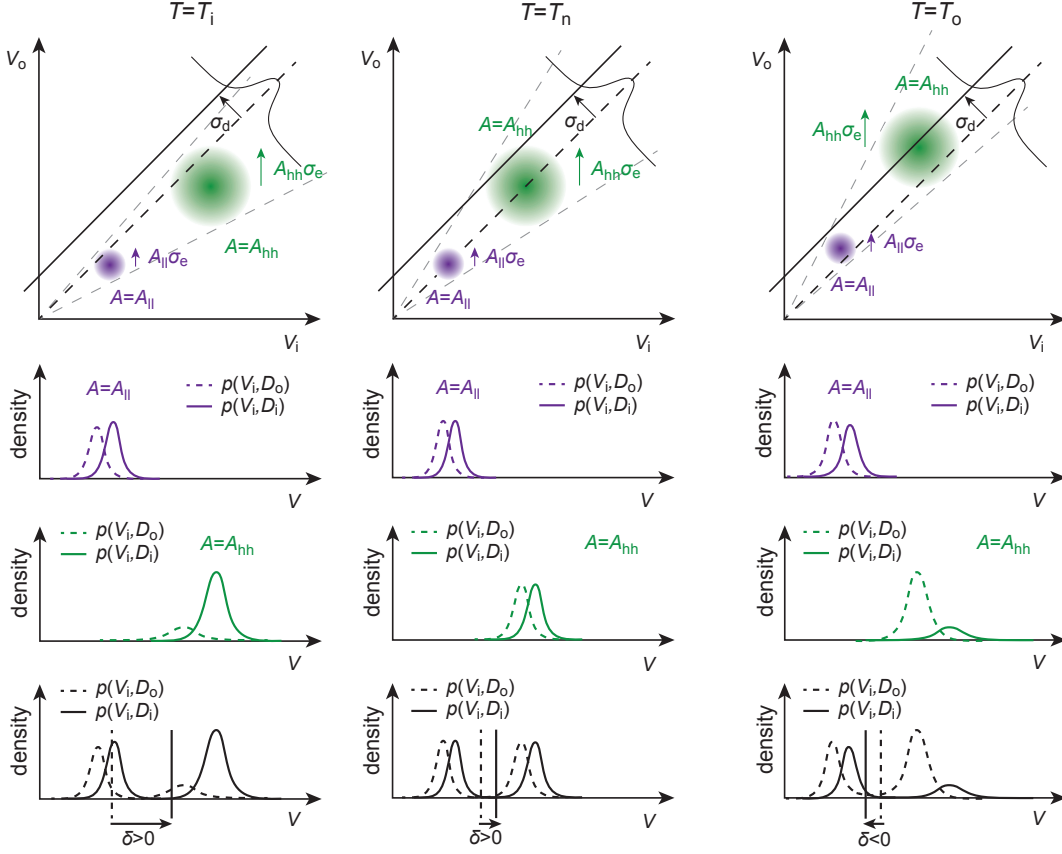


Figure 4: **Schematic for the multiplicative decision model with perfectly correlated fluctuating gain.**

We consider three stimulus conditions: “target-in”, i.e., $T = T_i$, “no-stimulus”, i.e., $T = T_n$, and “target-out”, i.e., $T = T_o$; and two equiprobable perfectly correlated attention states: “low gain”, i.e., $A = A_{||}$, and “high gain”, i.e., $A = A_{hh}$. The effect of varying attention states is to scale the distributions of voltage values (V_i, V_o) along the decision boundary in the two-dimensional phase space (top panels). Because of this scaling, the bias entailed by the presence of a stimulus is magnified in the high attention states A_{hh} . Specifically, in the high attention state, the voltage values (V_i, V_o) are more dispersed but stay further away from the decision boundary (solid straight line) on average. This is by contrast with the additive model of Fig. 3. As in Fig. 3, given an attention state, the recorded voltage data can be obtained by the same marginalization operation of the distribution of (V_i, V_o) restricted to each decision region of the two-dimensional phase space as in Fig. 1. Given knowledge of the attention state, the expected value of V_i is always larger in the “choice in” condition compared with the “choice out” condition. This is similar as for the additive model of Fig. 3. However, the key difference with Fig. 3 is that in the multiplicative model, there are much less misses in the high-attention state compared with the low attention state. As a result, when conditioned on the decision outcome alone and without knowledge of the attention state, the expected value of V_i can be less for “choice in” condition than for “choice out” when $T = T_o$, leading to a negative δ -metric. This is because when $T = T_o$, “choices in” are misses, which are more likely to happen in the low attention state, for which the gain variable A —thus the voltage V_i —is low.

For instance, one can check that for the “no-stimulus” condition, we have

$$\begin{aligned}
\delta(T_n) &= \sum_{\alpha, \beta \in \{l, h\}} \left(\mathbb{E}[V_i D_i | A_{\alpha\beta}, T_n] \frac{\mathbb{P}[A_{\alpha\beta}]}{\mathbb{P}[D_i | T_n]} - \mathbb{E}[V_i D_o | A_{\alpha\beta}, T_n] \frac{\mathbb{P}[A_{\alpha\beta}]}{\mathbb{P}[D_o | T_n]} \right), \\
&= 2 \sum_{\alpha, \beta \in \{l, h\}} (\mathbb{E}[V_i D_i | A_{\alpha\beta}, T_n] - \mathbb{E}[V_i D_o | A_{\alpha\beta}, T_n]) \mathbb{P}[A_{\alpha\beta}], \\
&= \sum_{\alpha, \beta \in \{l, h\}} \left(\frac{\mathbb{E}[V_i D_i | A_{\alpha\beta}, T_n]}{\mathbb{P}[D_i | A_{\alpha\beta}, T_n]} - \frac{\mathbb{E}[V_i D_o | A_{\alpha\beta}, T_n]}{\mathbb{P}[D_o | A_{\alpha\beta}, T_n]} \right) \mathbb{P}[A_{\alpha\beta}], \\
&= \sum_{\alpha, \beta \in \{l, h\}} \delta(T_n, A_{\alpha\beta}) \mathbb{P}[A_{\alpha\beta}],
\end{aligned}$$

so that the δ -metric is guaranteed to be positive as the average of the δ -metrics $\delta(T_n, A_{\alpha\beta}) > 0$ obtained for the case of a fixed, single attention state $A_{\alpha, \beta}$. The expressions for “stimulus in” and “stimulus out” are slightly more complicated as they necessitate evaluating the marginal miss probability m . Specifically, we have

$$\begin{aligned}
\delta(T_i) &= \sum_{\alpha, \beta \in \{l, h\}} \left(\mathbb{E}[V_i D_i | A_{\alpha\beta}, T_i] \frac{\mathbb{P}[A_{\alpha\beta}]}{\mathbb{P}[D_i | T_i]} - \mathbb{E}[V_i D_o | A_{\alpha\beta}, T_i] \frac{\mathbb{P}[A_{\alpha\beta}]}{\mathbb{P}[D_o | T_i]} \right), \\
&= \frac{1}{1-m} \sum_{\alpha, \beta \in \{l, h\}} \mathbb{E}[V_i D_i | A_{\alpha\beta}, T_i] \mathbb{P}[A_{\alpha\beta}] - \frac{1}{m} \sum_{\alpha, \beta \in \{l, h\}} \mathbb{E}[V_i D_o | A_{\alpha\beta}, T_i] \mathbb{P}[A_{\alpha\beta}], \\
\delta(T_o) &= \sum_{\alpha, \beta \in \{l, h\}} \left(\mathbb{E}[V_i D_i | A_{\alpha\beta}, T_o] \frac{\mathbb{P}[A_{\alpha\beta}]}{\mathbb{P}[D_i | T_o]} - \mathbb{E}[V_i D_o | A_{\alpha\beta}, T_o] \frac{\mathbb{P}[A_{\alpha\beta}]}{\mathbb{P}[D_o | T_o]} \right), \\
&= \frac{1}{m} \sum_{\alpha, \beta \in \{l, h\}} \mathbb{E}[V_i D_i | A_{\alpha\beta}, T_o] \mathbb{P}[A_{\alpha\beta}] - \frac{1}{1-m} \sum_{\alpha, \beta \in \{l, h\}} \mathbb{E}[V_i D_o | A_{\alpha\beta}, T_o] \mathbb{P}[A_{\alpha\beta}].
\end{aligned}$$

In the following, we use the above expressions to make quantitative predictions about the δ -metric. It is also not too hard to use these expressions to derive the heuristic arguments that explains why $\delta(T_o)$ can be negative for fluctuating gain shared across both hemispheres. For perfect correlation across hemisphere, we have $q = 0$ and for equiprobable low- and high-attention states we have $p = 1/2$, so the marginal hit probability is $m = (m_{ll} + m_{hh})/2$. This implies that

$$\delta(T_o) = \frac{1}{2m} (\mathbb{E}[V_i D_i | A_{ll}, T_o] + \mathbb{E}[V_i D_i | A_{hh}, T_o]) - \frac{1}{2(1-m)} (\mathbb{E}[V_i D_o | A_{ll}, T_o] + \mathbb{E}[V_i D_o | A_{hh}, T_o]).$$

In general, we have $m_{ll} = \mathbb{P}[D_o | A_{ll}, T_i] = \mathbb{P}[D_i | A_{ll}, T_o] < 1/2$ and $m_{hh} = \mathbb{P}[D_o | A_{hh}, T_i] = \mathbb{P}[D_i | A_{hh}, T_o] < 1/2$. We typically consider that decision bias entailed by the presence of a stimulus is small in the low-attention state so that $\mathbb{P}[D_o | A_{ll}, T] = \mathbb{P}[D_i | A_{ll}, T] \simeq 1/2 \simeq m_{ll}$, independent of the stimulus T . This assumption becomes exact in the limit of large downstream noise $\sigma_d \rightarrow \infty$ or vanishing gain amplification $A_{ll} \rightarrow 0$. By contrast, the probability $m_{h,h}$ becomes negligible for large enough amplification $A_{hh} \rightarrow \infty$. This involves that $\mathbb{E}[V_i D_i | A_{hh}, T_o]$ is also negligible in this limit, so that

$$\delta(T_o) = (\mathbb{E}[V_i | D_i, A_{ll}, T_o]) - \mathbb{E}[V_i | D_o, A_{ll}, T_o] - \mathbb{E}[V_o D_i | D_i, A_{hh}, T_o],$$

where the parenthesized term is the conditional δ -metric in the low attention state A_{ll} . This term becomes negligible for small enough gain $A_{ll} \rightarrow 0$ and large enough downstream noise. Therefore, the only meaningful contribution comes from the rightmost negative term, leading to an overall negative unconditioned δ -metric.

3.3 Parametric dependences of the δ -metric

In Fig 5, we assume perfect correlation across hemispheres ($q = 0$) and we represent the dependence of the δ -metrics $\delta(T_i)$, $\delta(T_n)$, and $\delta(T_o)$ on three parameters: the high-attention multiplicative gain a_{hh} in Fig 5A, the probability to be in the high-attention state $p_{hh} = 1 - p$ in Fig 5B, and the variance of the downstream noise σ_d in Fig 5C. This graphical illustration shows that the occurrence of negative $\delta(T_o)$ is due to having fluctuating attention states of large enough magnitudes and frequency.

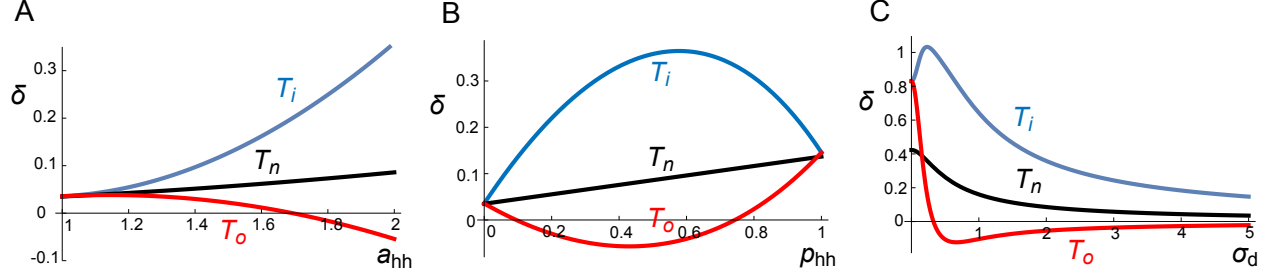


Figure 5: Parametric dependence of the δ -metric for the multiplicative model. **A.** Dependence on gain a_{hh} . For consistently low attention state $a_{hh} = a_{ll} = 1$, the δ -metric is almost constant at a low values, independent of the stimulus condition (actually, we have $\delta(T_n) < \delta(T_i) = \delta(T_o)$ but this distinction is negligible). When increasing the gain of the high attention state $a_{hh} > a_{ll} = 1$, the δ -metric separates for the three stimulus conditions, with negative values for “target out” when the gain is large enough. **B.** Dependence on high-attention state probability p_{hh} . The separation of the δ -metrics for the three stimulus conditions is due to gain fluctuations and is maximal for equiprobable low- and high-attention states. Observe that when the high-attention state is certain ($p_{hh} = 1$), δ -metrics revert to similar values. However, by contrast with the low attention state ($p_{hh} = 0$), the gain is large enough in the high attention state so that the δ -metrics for “target in” and “target out”, which are equal by symmetry, are noticeably larger than for the “no-stimulus” condition: $\delta(T_n) < \delta(T_i) = \delta(T_o)$. **C.** Dependence on the downstream noise σ_d . The limit of vanishing downstream noise $\sigma_d \rightarrow 0$ is akin to the limit of infinite gain a_{hh}, a_{ll} for which the δ -metrics for “target in” and “target out”, which are equal by symmetry, are significantly larger than for the “no-stimulus condition”. Increasing the downstream-noise variance σ_d generally reduces the δ -metrics leading to a regime for which the “target out” delta metric is negative for large enough values of σ_d . Parameters: $\mu_{V_i} = 2$, $\mu_{V_o} = 1$, $\sigma_{V_i} = 0.25$, $\sigma_{V_o} = 0.25$, $\sigma_d = 2$, $a_{ll} = 1$, $a_{hh} = 2$ and $p_{ll} = 1/2$. For all panels, both hemispheres are assumed perfectly correlated, i.e., $q = 0$.

In Fig 6, we relax the assumption of perfect correlations across hemispheres and explore under which conditions the δ -metric can become negative for the “target-out” condition. Our symmetric model for correlations involves two probability parameters p and q which are such that $1 - p - 2q \geq 0$. Thus the cross-hemisphere correlation $\rho(p, q)$ given via formula (1) is defined over a triangular region in the (p, q) -plane (see Fig. 6A). Fig 6B-D represent the δ -metrics $\delta(T_i)$, $\delta(T_n)$, and $\delta(T_o)$ as functions of the parameters p and q . The functional form of the δ -metric given above specifies the level sets of δ -metrics as quadrics. The “no-stimulus” case, for which the δ -metric is always positive, corresponds to the degenerate case when these level sets are affine functions. This graphical illustrations reveals that the occurrence of negative $\delta(T_o)$ is only consistent with $q \simeq 0$ and $p \simeq 1/2$, i.e., with a fluctuating attention state and with close to perfect correlation across hemispheres.

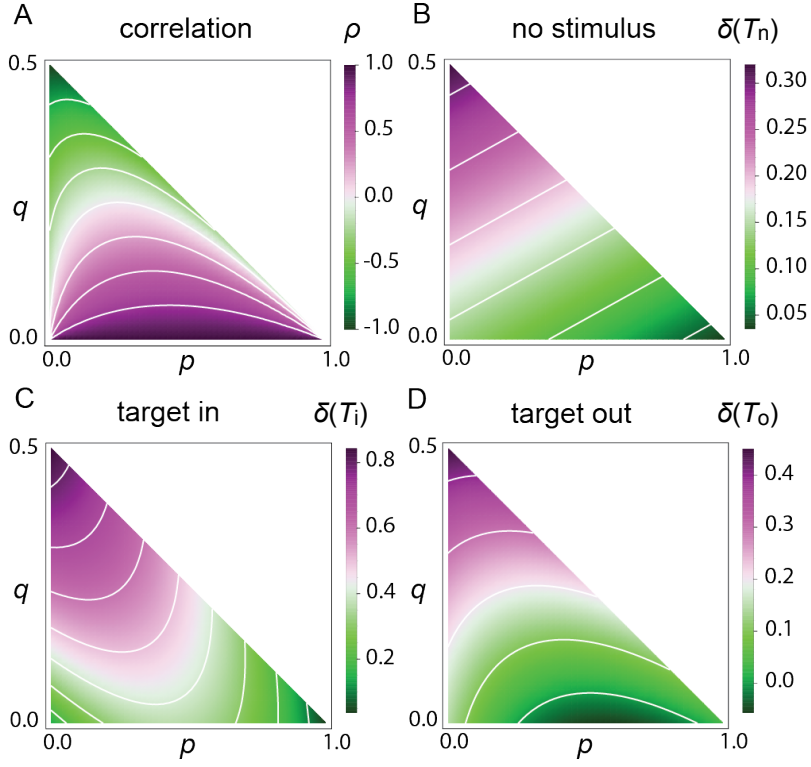


Figure 6: **Impact of correlations across hemispheres on δ -metric.** **A.** Correlation of the gain variables $\rho = \mathbb{C}[A_i, A_o] / \sqrt{\mathbb{V}[A_i] \mathbb{V}[A_o]}$ as a function of the probability parameters $p = \mathbb{P}[A_{ll}]$ and $q = \mathbb{P}[A_{lh}] = \mathbb{P}[A_{hl}]$. **B, C, D.** Delta metric $\delta(T_i)$ as a function of the same parameters p and q for the “no-stimulus” condition (**B**), and the “target-in” condition (**C**), and the “target-out” condition (**D**). For all stimulus conditions, the δ -metric increases with negative correlation across hemisphere, whereas the δ -metric generically satisfy $\delta(T_i) \geq \delta(T_n) \geq \delta(T_o)$. Negative δ -metric only occurs for $T = T_o$ and for strongly, fluctuating attention state, i.e., $q \simeq 0$ and $p \simeq 1/2$. Parameters: $\mu_{V_i} = 2$, $\mu_{V_o} = 1$, $\sigma_{V_i} = 0.25$, $\sigma_{V_o} = 0.25$, $\sigma_d = 2$, $a_{ll} = 1$, $a_{hh} = 2$.

3.4 Numerical results for d -prime metrics and choice probabilities

Our analysis was derived for a two-state gain model but generalizes straightforwardly to the case of a continuously varying attention state. To see this, we consider that the multiplicative gain A modeling attention is shared across hemispheres and follows a distribution with density

$$p(a) = \frac{1}{\sigma_a \sqrt{2\pi} a^3} \exp\left(-\frac{(1-a)^2}{2\sigma_a^2 a}\right) \quad \text{for } a \geq 0. \quad (2)$$

The above distribution is such that gain values A are positive with probability one, have unit expectation value, i.e. $\mathbb{E}[A] = 1$, and have variance σ_a^2 . Highly fluctuating attention states correspond to large values $\sigma_a^2 \gg 0$. One can extend the analytical results obtained for the δ -metric to continuously varying attention states by mere integration over the attention states. In Fig. 7A, we explore numerically the resulting expressions by varying σ_a , which quantifies the degree of gain fluctuation. As expected, we found that for the chosen set of parameters, $\delta(T_o)$ becomes negative for large enough σ_a .

In principle, one could further extend the analysis to the case of the d -prime metrics, a classical metric for choice discriminability. This would require, however, to compute conditional variances in addition to conditional expectations, which yields unwieldy expressions. Therefore, we only evaluate the d -prime metrics numerically in Fig. 7B, for the same set of parameters as in Fig. 7A. One can check that as expected, the d -prime metrics become negative in the “choice-out” condition for the same value σ_a as for the δ -metric, which is a scaled version of d -prime.

Our final goal is to check that our analysis remains qualitatively valid for another key choice-related metrics called choice probabilities. In our context, choice probabilities quantifies how well an ideal observer can predict the behavioral choice from reading out a neuron’s voltage distribution. Formally, given the stimulus condition T , this involves considering the choice-related conditional variables $U_i | T = V_i | (D_i, T)$ and $U_o | T = V_i | (D_o, T)$, where U_i refers to the voltage V_i conditioned on “choice-in” and U_o refers to the voltage V_i conditioned on “choice-out”. Then, the choice probabilities CP (T) is defined as the probability that U_i exceeds U_o :

$$\text{CP}(T) = \mathbb{P}[U_i \geq U_o | T] = \mathbb{P}[U_i | T \geq U_o | T].$$

In Fig. 7C, we compute numerically the choice probabilities associated to our multiplicative model for the same parameters as in Fig. 7A,B. Similar to Fig. 7A,B, one can check that for large enough σ_a , the choice probability CP (T_o) performs below chance level, a paradoxical effect akin to observing negative d -prime or δ -metric. Observe, however, that the paradoxical choice probabilities are observed for smaller values σ_a .

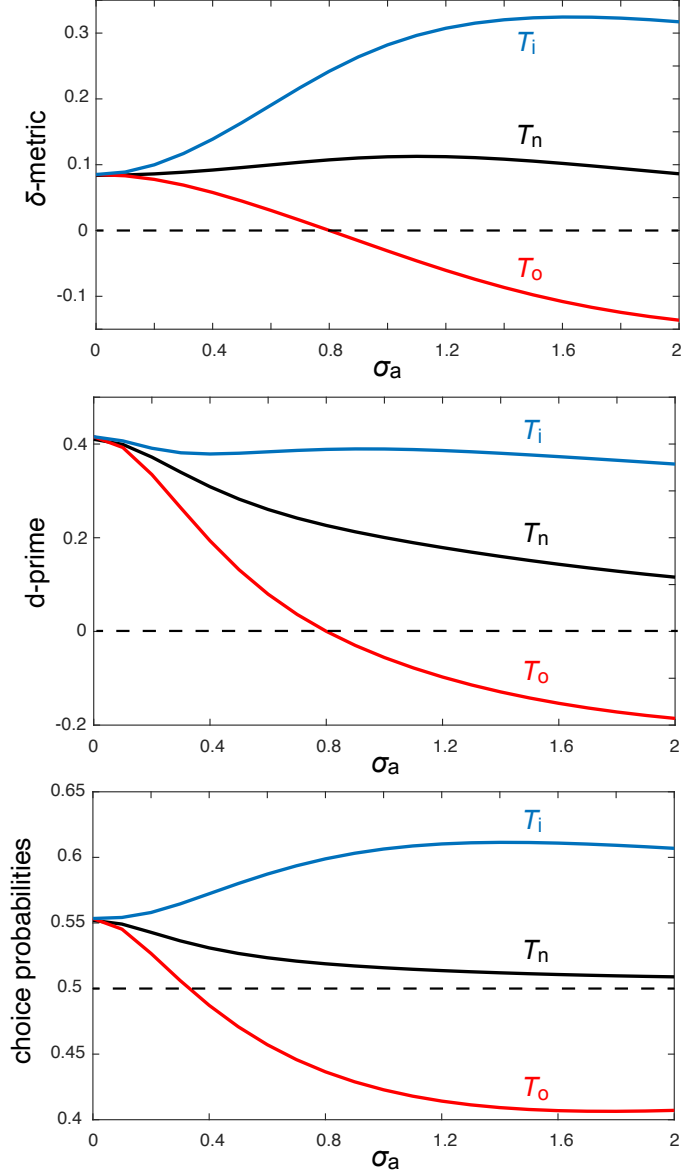


Figure 7: **δ -metric, d -prime metric, and choice probabilities.** The multiplicative gain is drawn from a distribution with density given by (2), which depends on a single variance parameter σ_a^2 . For all metrics, a paradoxical effect is observed in the “choice-out” condition for large enough fluctuations in the attention state: $\delta(T_o)$ and $d'(T_o)$ become negative, whereas CP (T_o) becomes smaller than chance level (0.5). Parameters: $\mu_{V_i} = 0.75$, $\mu_{V_o} = 0.5$, $\sigma_{V_i} = 0.1$, $\sigma_{V_o} = 0.1$, $\sigma_d = 1$.