

On the structure of neuronal population activity under fluctuations in attentional state

Alexander S. Ecker^{1,2,3,4}, George H. Denfield⁴, Matthias Bethge^{1,2,3,*}, and Andreas S. Tolias^{3,4,5,*}

¹Centre for Integrative Neuroscience and Institute for Theoretical Physics, University of Tübingen, Germany

²Max Planck Institute for Biological Cybernetics, Tübingen, Germany

³Bernstein Centre for Computational Neuroscience, Tübingen, Germany

⁴Department of Neuroscience, Baylor College of Medicine, Houston, TX, USA

⁵Department of Computational and Applied Mathematics, Rice University, Houston, TX, USA

*These authors contributed equally

April 21, 2015

Abstract

1 Attention is commonly thought to improve behavioral performance by increasing response gain and sup-
2 pressing shared variability in neuronal populations. However, both the focus and the strength of attention
3 are likely to vary from one experimental trial to the next, thereby inducing response variability unknown
4 to the experimenter. Here we study analytically how fluctuations in attentional state affect the structure
5 of population responses in a simple model of spatial and feature attention. In our model, attention acts
6 on the neural response exclusively by modulating each neuron's gain. Neurons are conditionally indepen-
7 dent given the stimulus and the attentional gain, and correlated activity arises only from trial-to-trial
8 fluctuations of the attentional state, which are unknown to the experimenter. We find that this simple
9 model can readily explain many aspects of neural response modulation under attention, such as increased
10 response gain, reduced individual and shared variability, increased correlations with firing rates, limited
11 range correlations, and differential correlations. We therefore suggest that attention may act primar-
12 ily by increasing response gain of individual neurons without affecting their correlation structure. The
13 experimentally observed reduction in correlations may instead result from reduced variability of the at-
14 tentional gain when a stimulus is attended. Moreover, we show that attentional gain fluctuations – even
15 if unknown to a downstream readout – do not impair the readout accuracy despite inducing limited-range
16 correlations.
17

18	Contents	
19	1 Introduction	3
20	2 Results	4
21	2.1 Fluctuations in spatial attention	4
22	2.2 Fluctuations of feature attention	6
23	2.3 Effect of attention-induced correlations on population coding	9
24	2.4 Identifying attentional fluctuations in experimental data	11
25	2.5 A new view on the reduction of shared variability under attention	13
26	3 Discussion	14
27	4 Appendix	15
28	4.1 Model setup	15
29	4.2 Effect of fluctuating gains on spike count statistics	16
30	4.3 Effect of fluctuations in attended feature on spike count statistics	16
31	4.4 Coding accuracy under fluctuations of spatial attention	17
32	4.5 Coding accuracy under fluctuations of feature attentional gain	18
33	4.6 Coding accuracy under fluctuations of attended feature	18

1 Introduction

Attention was traditionally thought of as acting by increasing response gain of a relevant population of neurons (Maunsell and Treue 2006; Reynolds and Chelazzi 2004). More recent studies found that attention also reduces pairwise correlations between neurons (Cohen and Maunsell 2009; Herrero et al. 2013; Mitchell et al. 2009). Based on a simple pooling model (Zohary et al. 1994) these authors argued that the effects of increased gain are dwarfed by the effects of reduced correlations and, therefore, attention is more appropriately viewed as shaping the noise distribution.

However, in an experiment the subject's state of attention can be controlled only indirectly and is bound to vary from one trial to the next. As a consequence, measuring neuronal variability or correlations under attention has a fundamental caveat: it is unclear to what extent the observed neuronal covariability reflects interesting aspects of information processing in the neuronal population or simply trial-to-trial fluctuations in the subject's state of attention, which is unknown to the experimenter. Despite ample evidence that attention fluctuates from trial to trial (Cohen and Maunsell 2010; Cohen and Maunsell 2011), the effects of such fluctuations on neuronal population activity have so far not been investigated.

Here we analyze a simple neural population model, where neurons with overlapping receptive fields encode the direction of motion of a stimulus. We assume that neurons produce spikes independently according to a Poisson process with rate λ_i and treat attention as a process that modulates the neurons' gain. The firing rates are given by

$$\lambda_i = g_i f_i(\theta), \quad (1)$$

where g_i is the attentional gain (a combination of spatial and feature attention) and $f_i(\theta)$ is the direction tuning curve of neuron i . We assume that there is always a stimulus in the neurons' receptive field, but this stimulus is not necessarily attended.

Crucially, in our model the subject's attentional state is not constant across trials, even within the same attentional condition. Thus, g_i is a random variable that varies from trial to trial, and its precise value is unknown to the experimenter. As a consequence, the correlations in g_i across neurons will induce correlations between the observed neural responses. In the following sections, we analyze this correlation structure in detail. In addition, we investigate the consequences of these correlations for reading out the direction of motion of the stimulus from the population response if the readout does not have access to the attentional state.

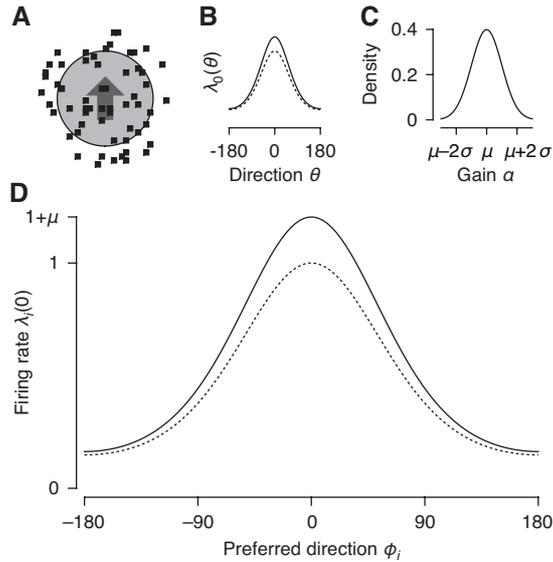


Figure 1. Model of spatial attention. **A.** Example stimulus. Neurons’ receptive fields are assumed to be at the same location (circle). **B.** Tuning curve under sensory stimulation (dashed) and with spatial attention directed to the stimulus in the receptive field (solid). **C.** Distribution of attentional gain (α). **D.** Population response of a homogeneous population of neurons under sensory stimulation (dashed) and with attention directed to the stimulus in the receptive fields (solid).

2 Results

2.1 Fluctuations in spatial attention

We first consider the simplest case of pure spatial attention and a common gain α for all neurons (Fig. 1):

$$\lambda_i = \alpha f_i(\theta), \quad (2)$$

where $\alpha > 0$ is the amount of spatial attention allocated to the stimulus in the neurons’ receptive field. We do not require any distributional assumptions on α , except for its mean $E[\alpha] = \mu$ and variance $\text{Var}[\alpha] = \sigma^2$ (Fig. 1C). Under this model, the average spike count of a neuron is given by

$$E[y_i|\theta] = \mu f_i(\theta). \quad (3)$$

By convention we refer to the case of $\mu = 1$ as the *sensory response*, which is the neural response to the stimulus in the absence of any attentional modulation. In experimental conditions where the stimulus is attended $\mu_a > 1$. When attention is directed towards a different stimulus $\mu_a \leq 1$ (depending on whether responses are suppressed relative to the sensory response under such conditions). Note that although we use homogeneous neural populations in the figures (all neurons have the same tuning curve up to a preferred direction ϕ_i , i.e. $f_i(\theta) = f(\theta - \phi_i)$), all results hold more generally for arbitrary tuning curves.

Because the attentional state fluctuates from trial to trial, the underlying firing rate also fluctuates. By applying the law of total variance we obtain the spike count variance (Fig. 2A):

$$\text{Var}[y_i|\theta] = \mu f_i(\theta) + \sigma^2 f_i^2(\theta). \quad (4)$$

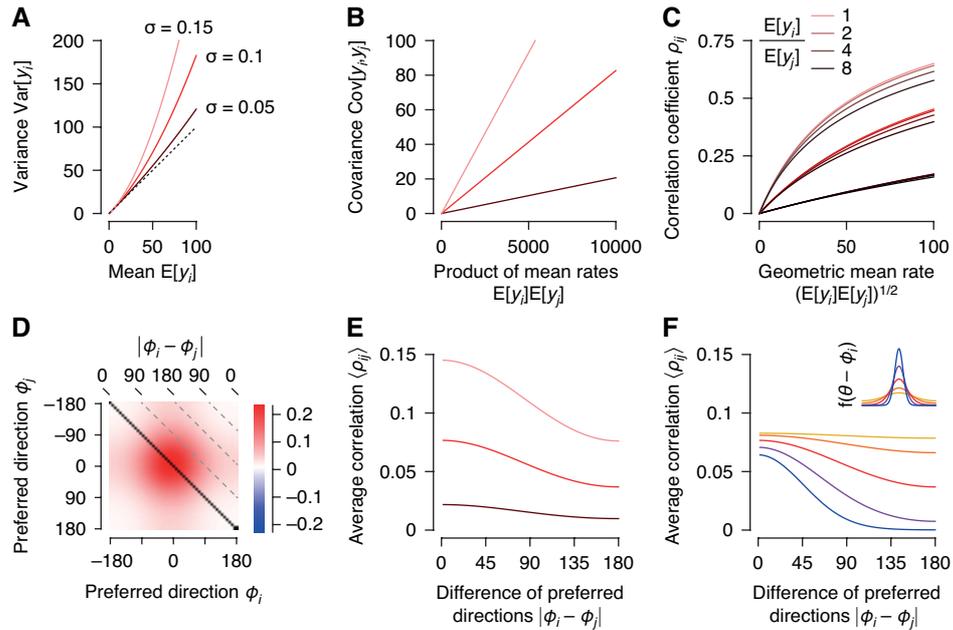


Figure 2. Effect of fluctuations in attentional state on spike count statistics. Solid lines: analytical solutions (Eqs. 3–7). Parameter values used here were $\mu = 0.1$, $\sigma \in \{0.05, 0.1, 0.15\}$ (dark to light red). **A.** Spike count variance as a function of mean spike count. Dashed line: identity (Poisson process). **B.** Covariance as a function of product of spike counts. **C.** Correlation coefficient as function of geometric mean firing rate. The three groups of lines correspond to different levels of σ as in the other panels. Darker colors within a group indicate increasing ratios $f_i(\theta)/f_j(\theta)$. **D.** Matrix of correlation coefficients for $\theta = 0^\circ$ and $\sigma = 0.1$. Tuning curves: $f_i(\theta) = \exp(\kappa \cos(\theta - \phi_i) + \epsilon)$, $\kappa = 2$, average firing rate across all θ : 10 spikes/s. **E.** Average correlation coefficient (over all directions of motion θ) as a function of difference of the preferred directions of the two neurons. Despite a common gain for all neurons, correlations decay with tuning difference. Parameters as in panel E. **F.** As in panel E, but for different tuning widths ($\kappa \in \{0.5, 1, 2, 4, 8\}$, shown in inset at the top). The decay of the correlations with the difference of the preferred directions is stronger for narrow tuning curves. Red line corresponds to panels D and E. Mean firing rate: 10 spikes/s for all tuning widths.

76 The first term is equal to the average spike count and results from the Poisson process assumption,
 77 while the second term is quadratic in the firing rate, which results from the multiplicative nature of the
 78 fluctuating gain α (Goris et al. 2014). Such an expanding mean-variance relation has been observed in
 79 many experimental studies (Britten et al. 1993; Dean 1981; Goris et al. 2014; Tolhurst et al. 1983). Note
 80 that if the attentional gain does not fluctuate, we recover the Poisson process.

81 Similar to the variances, we can compute the covariance between two neurons, which is given by the
 82 product of the firing rates and the variance of the attentional gain (Fig. 2B):

$$\text{Cov}[y_i, y_j | \theta] = \sigma^2 f_i(\theta) f_j(\theta) \quad i \neq j. \quad (5)$$

83 Recall that neurons are assumed to be conditionally independent given the attentional gain. Thus, any
 84 covariability arises exclusively from gain fluctuations. As a result, the covariance matrix (Fig. 2D) can
 85 be expressed as a diagonal matrix plus a rank-one matrix:

$$C = \mu \text{Diag}(\mathbf{f}) + \sigma^2 \mathbf{f} \mathbf{f}^T. \quad (6)$$

86 Note that the assumption of conditional independence could be relaxed without affecting any of the
87 major results qualitatively: the diagonal matrix in the equation above would simply be replaced by the
88 (non-diagonal) point process covariance matrix.

89 Experimental studies more typically quantify spike count correlations rather than covariances. We
90 therefore also calculated the correlation coefficient ρ_{ij} of two neurons (Fig. 2C):

$$\rho_{ij} = \sqrt{\frac{f_i f_j}{(\mu/\sigma^2 + f_i)(\mu/\sigma^2 + f_j)}} \quad (7)$$

91 The spike count correlations induced by a fluctuating attentional gain increase with firing rates $f_i(\theta)$.
92 This effect, which has also been observed in numerous experimental studies (Cohen and Maunsell 2009;
93 Ecker et al. 2014; Mitchell et al. 2009; Smith and Sommer 2013), arises because the independent (Poisson)
94 variability is linear in the firing rate, whereas the covariance induced by gain fluctuations is quadratic
95 and therefore dominates for large firing rates. Thus, correlations increase with the geometric mean firing
96 rate, but there is no simple one-to-one mapping between the two quantities (it also depends on the ratio
97 of the firing rates, Fig. 2C). The covariance, in contrast, is proportional to the product of the firing rates
98 with a constant of proportionality of σ^2 (Fig. 2B), suggesting that the latter might be more appropriate
99 to consider when analyzing experimental data.

100 In addition, the correlation structure induced by gain fluctuations is non-trivial even if all neurons
101 share the same gain (Fig. 2E, F; see also Ecker et al. (2014)). Due to the nonlinear shape of the tuning
102 function and the nonlinear way the neurons' tuning functions affect spike count correlations, the correla-
103 tions decrease with increased difference in two neurons' preferred directions (Fig. 2F). The slope of the
104 decay depends mainly on the dynamic range of the tuning curve. If neurons have a high baseline firing
105 rate compared to their peak firing rate, correlations decrease only marginally with preferred direction.
106 In contrast, sharply tuned neurons with close to zero baseline firing rates exhibit strong limited-range
107 structure.

108 This limited-range correlation structure has been observed in numerous experimental studies (Bair
109 et al. 2001; Cohen and Maunsell 2009; Ecker et al. 2010; Smith and Kohn 2008; Zohary et al. 1994) and
110 has been hypothesized to reflect shared input among similarly tuned neurons. However, our simple model
111 shows that these seemingly structured correlations can arise from a very simple, non-specific mechanism:
112 a common fluctuating gain that drives all neurons equally, irrespective of their tuning properties.

113 2.2 Fluctuations of feature attention

114 Feature attention is different from spatial attention in that the sign of the gain modulation depends on
115 the similarity of the attended direction to the neuron's preferred direction of motion (Fig. 3). Following
116 the feature-similarity gain model (Treue and Martinez-Trujillo 1999), we model feature attention by

$$\lambda_i = (1 + \beta h_i(\psi)) f_i(\theta), \quad (8)$$

117 where β is the *feature gain* that controls how strongly the feature ψ (in this case direction of motion) is
118 attended on the given trial and $h_i(\psi)$ is the *gain profile* (Fig. 3B) that determines the sign and relative
119 strength of modulation for each neuron depending on the similarity of its preferred direction ϕ_i to the
120 attended direction ψ . We assume that $h_i(\psi)$ most strongly enhances neurons with preferred directions
121 equal to the attended direction and suppresses those with opposite preferred directions (Fig. 3B).

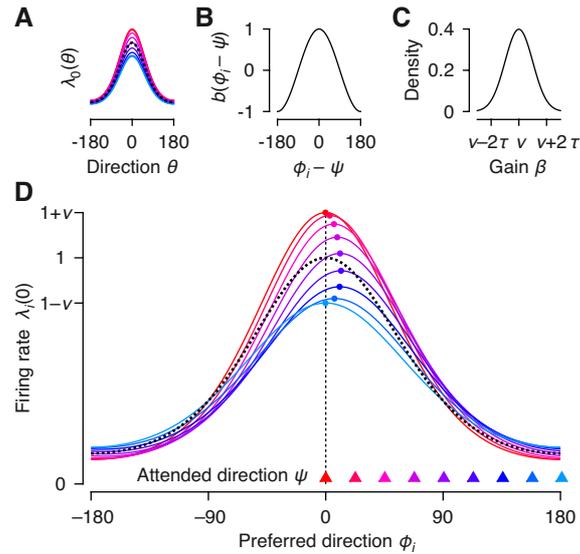


Figure 3. Model of feature attention. **A.** Tuning curve of a single neuron under sensory stimulation (black dotted) and with feature attention directed to different directions ranging from preferred (red) to null (blue). Note that the entire tuning curve of the neuron is gain-modulated and the modulation does not depend on the stimulus θ . **B.** The gain of a neuron depends on which direction of motion ψ is attended relative to the neuron’s preferred direction ϕ_i . **C.** Distribution of gain (β) fluctuations. A Gaussian is shown for illustration purposes; the analysis holds for any distribution with $E[\beta] = \nu$ and $\text{Var}[\beta] = \tau^2$. **D.** Population response of a homogeneous population of neurons under sensory stimulation (black dotted) and with attention directed to different directions of motion ranging from 0° (red) to 180° (blue). The stimulus is $\theta = 0$. The curves show the average response of the neurons as a function of their preferred direction. Attending to a direction of motion biases the population response towards this attended stimulus. While each neuron’s tuning curve is gain-modulated as a whole (panel A), the population response is no longer equal to the individual neurons’ tuning curves, but instead sharpened/broadened and its peak is moved.

122 Because feature attention both increases and decreases different neurons’ gain depending on their
 123 preferred direction relative to the attended direction of motion, it biases the population response towards
 124 the attended direction (Fig. 3D). Thus, unlike in the case of spatial attention the shape of the population
 125 response is no longer identical to that of the individual neuron’s tuning curve. We start by assuming that
 126 the subject always attends the same direction (i. e. ψ is constant) and consider the effect of fluctuations in
 127 the strength of attention, that is the gain β . We will come back to fluctuations in the attended direction
 128 below.

Similar to spatial attention, fluctuations in feature attention lead to overdispersion of the spike counts relative to a Poisson process (because rate variability is added).

$$E[y_i|\theta, \psi] = (1 + \nu h_i(\psi))f_i(\theta) \quad (9)$$

$$\text{Var}[y_i|\theta, \psi] = (1 + \nu h_i(\psi))f_i(\theta) + \tau^2 h_i^2(\psi) f_i^2(\theta), \quad (10)$$

129 where $\nu = E[\beta]$ and $\tau^2 = \text{Var}[\beta]$ are the mean and the variance of the feature attention gain, respectively.
 130 The degree of overdispersion not only increases with the neuron’s firing rate, but also depends on the
 131 neuron’s preferred direction relative to the attended direction (Fig. 4A). Interestingly, spike counts are

132 more overdispersed at the null direction than at the preferred direction (Fig. 4A: compare blue vs. black
133 and green vs. yellow). The Fano factor (variance/mean) is given by

$$F[y_i|\theta, \psi] = 1 + \frac{\tau^2 h_i^2(\psi)}{(1 + \nu h_i(\psi))^2} E[y_i], \quad (11)$$

134 which is higher when h_i is negative than when it is positive. Neurons with preferred directions orthogonal
135 to the attended direction are not overdispersed since $h_i = 0$.

136 As feature attention induces both increases as well as decreases in neuronal gain, the induced corre-
137 lation structure is different from that induced by spatial attention. For the covariances, we obtain

$$\text{Cov}[y_i, y_j|\theta, \psi] = \tau^2 h_i(\psi) h_j(\psi) f_i(\theta) f_j(\theta) \quad i \neq j \quad (12)$$

138 The sign of the covariance is determined by the product of h_i and h_j , which depends on the attended
139 direction and the preferred directions of the two neurons (Fig. 4B). For two neurons with identical
140 preferred directions, the covariance is always positive while for two neurons with orthogonal preferred
141 directions it is always negative. For any pair of neurons in between, it can be both positive and negative,
142 depending on the stimulus (Fig. 4B). Again, the covariance matrix can be written as diagonal plus rank
143 one:

$$C = F + \tau^2 \mathbf{u}\mathbf{u}^T, \quad (13)$$

144 where $F_{ii} = (1 + \nu h_i(\psi)) f_i(\theta)$ and $u_i = h_i(\psi) f_i(\theta)$.

145 As for spatial attention, averaging correlations over multiple stimulus conditions to represent the cor-
146 relation structure as a function of the neurons' tuning similarity misses much of the underlying structure
147 (Fig. 4C): spike count correlations are positively correlated with tuning similarity (Fig. 4D), but the
148 stimulus dependence (Fig. 4C) is again ignored. As before, the exact shape of the decay depends on
149 the tuning width: for narrow tuning curves, neurons with opposite preferred directions are only weakly
150 anti-correlated, whereas for broad tuning curves, those neurons are strongly anti-correlated (Fig. 4D, blue
151 to red lines).

So far we have assumed that the attended direction of motion is constant and only the strength of
attention fluctuates from trial to trial. Now we turn to the case where the attended direction fluctuates
from trial to trial. We assume that, on average, the subject attends the correct direction, i. e. $E[\psi] = \theta$,
but with some variance $\text{Var}[\psi] = q^2$. We further assume the gain β is constant. In this case, means and
covariances of the observed spike counts are given by

$$E[y_i|\theta, \beta] = (1 + \beta h_i) f_i \quad (14)$$

$$\text{Cov}[y_i, y_j|\theta, \beta] = \delta_{ij} (1 + \beta h_i) f_i + q^2 \beta^2 h'_i h'_j f_i f_j, \quad (15)$$

152 where $h'_i = \frac{d}{d\psi} h_i$ and we have abbreviated $h_i \equiv h_i(\theta)$ and $f_i \equiv f_i(\theta)$. As before, we can write the
153 covariance matrix as diagonal plus rank one:

$$C = F + q^2 \mathbf{v}\mathbf{v}^T, \quad (16)$$

154 where $F_{ii} = (1 + \beta h_i) f_i$ and $v_i = \beta h'_i f_i$. This pattern of correlations (Fig. 5) differs from those observed
155 before for gain fluctuations in an important way: the sign of the correlation between two neurons depends
156 only on whether their preferred directions are on the same side (both clockwise or counter-clockwise) of

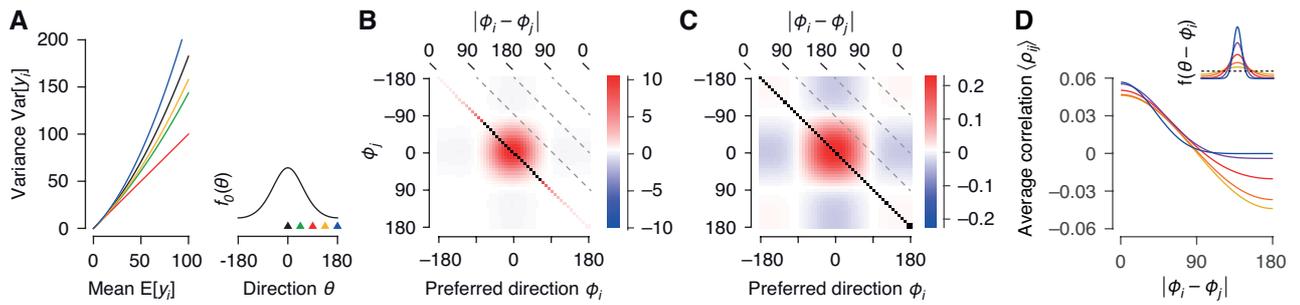


Figure 4. Effect of fluctuations in the feature attention gain on spike count statistics. Parameters here are: $\psi = 0$, $\nu = 0.1$, $\tau^2 = 0.01$. **A.** Spike count variance as a function of mean spike count. Colors indicate different attended directions relative to the neurons’ preferred direction ($\phi_i - \psi$; illustrated by colored triangles in inset on the bottom right). **B.** Covariance matrix for stimulus $\theta = 0$. Neurons are ordered by preferred directions. Mean firing rate across the population: 20 spikes/s. **C.** As panel B, but the correlation coefficient matrix is shown. **D.** Dependence of spike count correlations on tuning similarity (difference of preferred directions). Fluctuations in feature attention induce limited range correlations irrespective of the shape of the tuning curve. The higher the baseline firing rate the stronger the negative correlations for neurons with opposite preferred directions. Inset: different tuning widths used.

157 the stimulus direction or on different sides. As we will show more formally in the next section, this pattern
 158 of correlations is known as *differential correlations* (Moreno-Bote et al. 2014). Again, when plotted as a
 159 function of the difference of two neurons’ preferred directions, the correlations exhibit the typical limited-
 160 range structure (Fig. 5C), except for very narrow tuning curves, where the correlations are minimal
 161 around pairs with orthogonal preferred directions (Fig. 5C, blue lines). Also note that these correlations
 162 are substantially weaker than those induced by gain fluctuations (Figs. 2, 4), despite a relatively wide
 163 distribution of attended directions (SD: 10°).

164 2.3 Effect of attention-induced correlations on population coding

165 How interneuronal correlations affect the representational accuracy of neuronal populations has been a
 166 matter of immense interest (and debate) over the last years. Thus, we want to briefly consider how
 167 correlations induced by attentional fluctuations affect the coding accuracy of a population code.

168 Before doing so we need to make a choice: does the downstream readout have access to the state of
 169 attention or not? If it does, the picture is fairly simple: attentional fluctuations do not affect the readout
 170 accuracy, since the attentional state can be accounted for and there is no additional noise compared with
 171 a scenario without attentional fluctuations. The only downside is a potentially more complex readout. In
 172 contrast, if we assume that the readout does not have access to the attentional state, the situation becomes
 173 more interesting. In this case the attentional fluctuations act like additional (internally generated) noise,
 174 which could impair the readout. In the following we consider this latter scenario.

175 To quantify the accuracy of a population code, we use the Fisher information (Kay 1993) with respect
 176 to direction of motion. The Fisher information is useful because it quantifies the amount of information
 177 in a population of neurons without assuming a specific decoder. For a population of independent neurons,
 178 the Fisher information is linear in the number of neurons.

179 We start by considering spatial attention. Since the gain is the same for all neurons, gain fluctuations
 180 should not affect the coding accuracy of the population with respect to the direction of the stimulus,

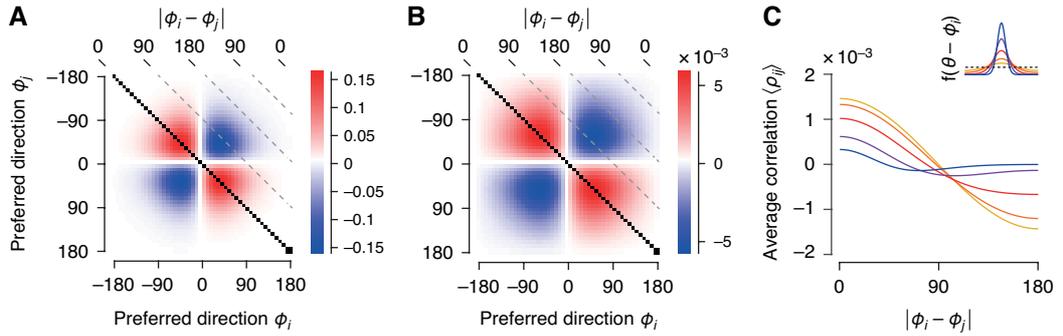


Figure 5. Effect of fluctuations in the attended direction on correlation structure. Parameters here are: $E[\psi] = 0$, $q = 10^\circ$, $\beta = 0.1$, $\theta = 0$, mean firing rate across the population: 20 spikes/s. **A.** Covariance matrix. Neurons are ordered by preferred directions. **B.** As panel A, but the correlation coefficient matrix is shown. **C.** Dependence of spike count correlations on tuning similarity (difference of preferred directions). Fluctuations in the attended direction induce limited range correlations, whose shape depends on the width of the tuning curves. Inset: different tuning widths used.

181 which is encoded in the differential activation pattern of the neurons. This is indeed the case. The Fisher
 182 information of a population of Poisson neurons whose firing rates are modulated by a common gain with
 183 mean μ is given by

$$J = \mu \sum_{i=1}^N \frac{f_i'(\theta)^2}{f_i(\theta)} - O(1) \approx J_0. \quad (17)$$

184 Thus, unobserved gain fluctuations reduce the information in the population only by a constant term (for
 185 derivation see Appendix). For reasonably large populations (e.g. $N > 100$) this term can be neglected
 186 and the information is approximately equal to that of an independent population (J_0). This result can
 187 be understood intuitively by considering the structure of the covariance matrix (Eq. 5): the dominant
 188 eigenvector points in the direction of the tuning function \mathbf{f} , which is orthogonal to changes in the stimulus,
 189 \mathbf{f}' . Therefore, gain fluctuations do not impair the readout of the direction of motion.

190 The same result holds for fluctuations in the feature attention gain, so long as the attended direction
 191 matches the one shown and does not fluctuate from trial to trial. A fluctuating gain sharpens and
 192 broadens the population hill from trial to trial, but leaves its peak unchanged. Again, the dominant
 193 eigenvector ($u_i = h_i f_i$, Eq. 12) points in a direction that is orthogonal to changes in the stimulus (details
 194 see Appendix).

195 The situation changes if the focus of attention (i.e. the attended direction) fluctuates from trial to
 196 trial or the attended direction does not match the one shown: since feature attention biases the population
 197 response towards the attended direction, such attentional fluctuations have the same effect as noise on
 198 the input [*differential correlations*, (Moreno-Bote et al. 2014)]. To illustrate this finding, we switch to a
 199 slightly modified and more specific response model than above. Assuming $f_i(\theta) = \exp(\kappa \cos(\theta - \phi_i))$ and
 200 $h_i = \cos(\psi - \phi_i)$, and noting that $(1 + \beta h_i) \approx \exp(\beta h_i)$, we can write the log-firing rate as

$$\log \lambda_i = \beta \cos(\psi - \phi_i) + \kappa \cos(\theta - \phi_i). \quad (18)$$

201 We can combine the two cosine terms and obtain:

$$\log \lambda_i = \gamma \cos(\theta + \Delta\theta - \phi_i) \quad (19)$$

where

$$\gamma = \sqrt{\kappa^2 + \beta^2 + 2\kappa\beta \cos(\psi - \theta)} \quad (20)$$

$$\Delta\theta = \arccos\left(\frac{\gamma^2 + \kappa^2 - \beta^2}{2\gamma\kappa}\right). \quad (21)$$

202 Thus, feature attention biases the population response away from the stimulus direction θ towards the
203 attended direction ψ . The magnitude of the bias $\Delta\theta$ depends on both the strength of feature attention
204 β and the attended direction ψ . Consequently, if $\psi \neq \theta$ fluctuations in either the attended feature or the
205 degree of feature attention have the same effect on the population response as variance of the stimulus
206 direction that is shown, i. e. they induce differential correlations. This result can also be understood by
207 considering the structure of the covariance matrix (Eq. 16): the dominant eigenvector $v_i = h'_i f_i$ points in
208 the same direction as changes in the stimulus, \mathbf{f}' . We can therefore approximate the Fisher information
209 by (see Moreno-Bote et al. 2014)

$$J \approx \frac{J_0}{1 + \varepsilon J_0} \rightarrow \frac{1}{\varepsilon}, \quad (22)$$

210 where J_0 is again the information in an independent population and $\varepsilon = \text{Var}[\Delta\theta]$ depends on both the
211 distribution of attended directions and the variance of the gain. In this case, the information in the
212 population saturates at a finite value $1/\varepsilon$ that depends only on the distribution of the attention signal
213 and can be substantially lower than the limit imposed by the information in feedforward signal (see also
214 Discussion). When the subject attends the correct direction on average (i. e. $\text{E}[\psi] = \theta$) and the variance
215 of the attended direction ($\text{Var}[\psi]$) is small, we find

$$J \rightarrow \frac{(\kappa/\beta + 1)^2}{\text{Var}[\psi]}. \quad (23)$$

216 Thus, the saturation level depends on the strength (β) of attention relative to the tuning width (κ) and
217 the variance in the attended direction.

218 2.4 Identifying attentional fluctuations in experimental data

219 We saw above that fluctuations in attentional state can introduce interesting patterns of correlations in
220 neural activity, all of which are roughly consistent with the published literature on attention. However,
221 as long as one considers only single neurons and pairwise statistics, any result can be consistent with
222 many hypotheses. For instance, attentional fluctuations induce correlations that depend on firing rates
223 (Fig. 2C), but the same result is also predicted by the thresholding nonlinearity of neurons (Rocha et al.
224 2007) and therefore need not result from attentional fluctuations. Similarly, all types of attentional fluc-
225 tuations considered above lead to correlations that decrease with the difference of two neurons' preferred
226 directions (*limited range correlations*, Figs. 2E, 4D, 5C), but this correlation structure can also arise from
227 shared sensory noise (Shadlen and Newsome 1998).

228 So how would one go about identifying attentional fluctuations in experimental data? Clearly, one
229 has to consider the response patterns of simultaneously recorded populations of neurons rather than just
230 pairwise correlations. In the following, we discuss some predictions our model makes for the structure of
231 the neural population response.

232 A first approach suggested by our analyses above: we showed that in all cases we analyzed the covari-
233 ance matrix induced by attentional fluctuations is diagonal plus rank one. Thus, attentional fluctuations

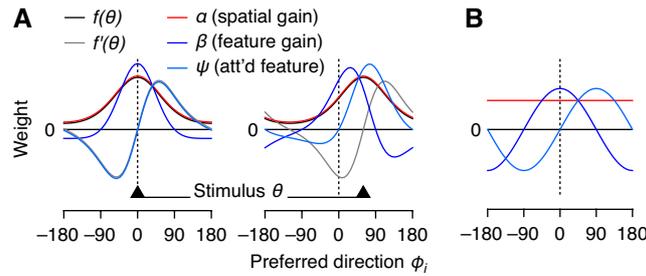


Figure 6. Identifying attentional fluctuations from variability in neuronal population activity. **A.** The subspace identified by Factor Analysis depends on the stimulus direction. Black triangles: stimulus direction (left: $\theta = 0^\circ$, right: $\theta = 60^\circ$). Solid lines: basis functions corresponding to fluctuations in spatial attention gain (red), feature attention gain (dark blue), and attended direction (light blue); population tuning curve (black) and its derivative (gray). Horizontal dashed line: (average) attended direction. **B.** Principal components identified by Exponential Family PCA are independent of the stimulus since the log-link turns a multiplicative modulation into an additive offset. Colors as in panel A.

234 are restricted to a low-dimensional subspace that could be identified from simultaneously recorded neu-
 235 rons by *Factor Analysis*. However, the disadvantage of this approach is that the low-dimensional subspace
 236 depends on the stimulus in a non-trivial way (Fig. 6A for $\theta = 0^\circ$ [left] and $\theta = 60^\circ$ [right]; see also Eqs. 5,
 237 12, 16). This stimulus dependence precludes pooling of data over multiple stimulus conditions. Moreover,
 238 if the attended direction does not match the stimulus direction, the major axes of variability do not peak
 239 at either direction, but somewhere in between (Fig. 6A, blue lines in the right panel, where $\psi = 0^\circ$ and
 240 $\theta = 60^\circ$). Thus, it is non-trivial to recover the quantities of interest for the experimenter – the attended
 241 feature (direction) and the degree of attention allocated (the gain).

242 A model that could directly extract attentional gains (spatial and feature gain) and the attended
 243 feature would be desirable. Fortunately, all three can be inferred from population activity in a straight-
 244 forward manner using methods such as *Exponential Family Principal Component Analysis (E-PCA)*
 245 (Collins et al. 2001; Mohamed et al. 2009) or *Poisson Linear Dynamical Systems (PLDS)* (Buesing et al.
 246 2012; Macke et al. 2011). Similar to above (Eq. 18), we assume $f_i(\theta) = \exp(\kappa_i \cos(\theta - \phi_i) + \epsilon_i)$ and
 247 $h_i = \cos(\psi - \phi_i)$ and write the log-firing rate as

$$\log \lambda_i = \alpha + \beta \cos(\psi - \phi_i) + \kappa_i \cos(\theta - \phi_i) + \epsilon_i, \quad (24)$$

248 which can be rewritten as a linear function of the attentional state and the stimulus:

$$\log \lambda_i = a + \mathbf{k}_i^T \mathbf{b} + \kappa_i \mathbf{k}_i^T \mathbf{x} + \epsilon_i, \quad (25)$$

249 where a and $\mathbf{b} = \beta \cdot [\cos \psi, \sin \psi]^T$ represent the state of spatial and feature attention, respectively,
 250 $\mathbf{x} = [\cos \theta, \sin \theta]^T$ is the stimulus, $\mathbf{k}_i = [\cos \phi_i, \sin \phi_i]^T$ is the neuron’s preferred direction, κ_i the
 251 (inverse) tuning width, and ϵ_i controls the mean firing rate. This model is a *Generalized Linear Model*
 252 (*GLM*) with Poisson observations and $\log(x)$ as the link function. Thus, E-PCA or PLDS will recover
 253 the subspace corresponding to fluctuations in attentional state $\{a, \mathbf{b}\}$. This subspace is spanned by
 254 $\mathbf{u}_i = [1, \cos \phi_i, \sin \phi_i]$ and independent of the stimulus (see Fig. 6B). The attentional gains are a and
 255 $\beta = \|\mathbf{b}\|$, while the attended direction is $\psi = \angle \mathbf{b}$.

256 2.5 A new view on the reduction of shared variability under attention

257 There is ample experimental evidence that attention fluctuates from trial to trial (Cohen and Maunsell
258 2010; Cohen and Maunsell 2011), and we showed in the previous sections that such fluctuations induce
259 patterns of (correlated) variability that are highly consistent with the reported data on attention (Cohen
260 and Maunsell 2009; Herrero et al. 2013; Mitchell et al. 2009). Interestingly, in our model, both the
261 magnitude of overdispersion in single neurons' spike counts and the average level of correlations are
262 determined by the variance of the attentional gain ($\sigma^2 = \text{Var}[g]$), but not by its average modulation
263 ($\mu = \text{E}[g]$). This observation suggests that the average attentional modulation (μ) between an attended
264 and an unattended condition (which can be reliably measured based on average responses) does not
265 predict the level of correlations in either condition, since the latter is controlled by an independent
266 variable (σ^2). Indeed, this is one of the central experimental findings: directing spatial attention to a
267 certain location *increases* the average responses of neurons whose receptive fields represent this location,
268 but reduces independent and shared variability among those neurons (Cohen and Maunsell 2009; Herrero
269 et al. 2013; Mitchell et al. 2009). Thus, if our model is correct, then the data suggest that attention not
270 only increases response gain, but also reduces the trial-to-trial variability of the gain.

271 This view of attention has important implications for the role of interneuronal correlations under
272 attention. Recent studies (Cohen and Maunsell 2009; Mitchell et al. 2009) have argued that spatial
273 attention improves behavioral performance primarily by reducing correlations. However, as we showed
274 above, fluctuations of spatial attention do not affect the representational accuracy of the neuronal pop-
275 ulation. Therefore, under our model the experimentally observed reduction in correlations is irrelevant
276 when reading out a neuronal population. The only difference that matters is the increase in gain.

3 Discussion

We find that a simple model of neuronal responses can account for a range of empirically observed phenomena relating attention, neuronal variability and coding properties of neuronal populations. Our model unites two central findings in the literature on attention, that attention acts as a multiplicative gain factor on neuronal responses (Maunsell and Treue 2006) and that attention fluctuates from trial-to-trial (Cohen and Maunsell 2010). The importance of the combined effects of these observations has not previously been fully appreciated, as we show that such a model is sufficient to account for super-Poisson variability (see also Ecker and Tolias 2014; Goris et al. 2014) as well as a variety of pairwise correlation structures, most notably the “limited-range” structure and “differential correlations” (Abbott and Dayan 1999; Ecker et al. 2010; Moreno-Bote et al. 2014; Smith and Kohn 2008).

Our results argue that it is likely that a large fraction of variability in the neuronal response can be attributed to fluctuations in behaviorally relevant, internally-generated signals, such as attention, rather than shared noise (Ecker and Tolias 2014; Ecker et al. 2010, 2014; Goris et al. 2014; Nienborg and Cumming 2009). This view suggests the hypothesis that correlations that arise from such fluctuating signals generally should not impair coding of sensory information. We find that this assertion is true for the case of fluctuations in the magnitude of the gain. The Fisher information of our model population of neurons is not limited by fluctuations in the strength of attentional gain (i.e., is independent of the variance of the gain term), despite those fluctuations generating a “limited-range” correlation structure typically thought to impair coding.

However, theoretical work has shown that the effect of different patterns of correlations on the coding of sensory information is nuanced and can depend greatly on specific assumptions that are made regarding a variety of neuronal properties, such as the shapes of tuning curves in the population, subtle details of the assumed correlation structure, or different readouts (Abbott and Dayan 1999; Ecker et al. 2011; Josić et al. 2009; Shamir and Sompolinsky 2006; Sompolinsky et al. 2001; Wilke and Eurich 2002). The recent work of Moreno-Bote et al. (2014) has helped to clarify the problem of when and what types of correlation structures are detrimental to coding with their description of “differential correlations,” a specific pattern of correlation proportional to the product of the derivative of the tuning curves that leads to information saturation. Our model generates this pattern of correlated variability when the fluctuations in attention occur around a specific feature rather than a specific gain value. Thus, it is noteworthy that a model only slightly more complicated than typical Poisson spiking models can generate the diversity of correlation structures noted in the experimental and theoretical literature as being important for population coding.

In addition to offering a parsimonious account of neuronal variability and co-variability, our model has implications for how we should interpret the effect of attention as it relates to improvements in perceptual performance. Chiefly, if the reduction of correlations observed under attention is indeed due to a reduction of gain fluctuations – as our model would suggest – the reduction of correlations is irrelevant with respect to the coding accuracy of the population and cannot be the mechanism improving behavioral performance as suggested by recent experimental studies (Cohen and Maunsell 2009; Herrero et al. 2013; Mitchell et al. 2009).

Our model leads to a second interesting observation: It is likely that not only the attentional gain fluctuates from trial to trial, but also the attended feature itself. Such fluctuations introduce differential correlations, which indeed impair the readout (unless it has exact access to the attended feature). Thus, the attentional mechanism itself places a limit on how accurately a stimulus can be represented by a sensory population, and this limit can at least in principle be substantially lower than the amount of

320 sensory information entering the brain through the eye. This insight may trigger the question: why, then,
321 should there be an attentional mechanism in the first place? There are a number of possible answers to
322 this question.

323 First, we can think of attention as a prior. Using prior information to bias an estimate towards
324 more likely solutions will on average improve the estimate. In situations where the stimulus is noisy and
325 decisions have to be made fast, such a bias is most beneficial and outweighs the small extra noise added
326 due to variability in the prior. Conversely, in situations where there is lots of sensory evidence, the full
327 information content present in the eye is rarely necessary in real-world situations, and, therefore, the
328 noise added due to attentional fluctuations does not matter either.

329 Second, it should be noted that for change-detection paradigms that are typically employed in atten-
330 tion experiments, the estimation framework that asks how well a stimulus value can be reconstructed
331 (e.g. Fisher information) is not quite appropriate. In such tasks the subject never judges the *absolute*
332 direction (or any other feature) of the stimulus, but instead has to detect a small change, that is the
333 difference between two subsequent stimuli. In this case any errors introduced due to fluctuations in the
334 attended direction cancel out, since they affect both stimuli roughly equally, at least so long as attentional
335 fluctuations occur at a timescale that is slow enough, such that the attentional state is approximately
336 the same for both the pre- and post-change stimulus.

337 4 Appendix

338 4.1 Model setup

339 We model a population of direction-selective neurons with overlapping receptive fields and a diverse range
340 of preferred directions ϕ_i . We use a simple model of spatial and feature attention, where a neuron's firing
341 rate λ_i is the product of an attentional gain $g_i(\psi)$ and a tuning function $f_i(\theta)$:

$$\lambda_i(\theta, \psi) = g_i(\psi)f_i(\theta) \quad (26)$$

342 Here, ψ is the attended direction of motion and θ the direction of the stimulus that is shown. Neurons
343 are assumed to be conditionally independent given the firing rate λ_i (i.e. no noise correlations). The
344 attentional gain depends on whether attention is directed to the location of the neurons' receptive fields
345 and on the attended direction of motion. For spatial attention, we use $g_i = \alpha$, which is the same
346 for all neurons, since they all have overlapping receptive fields. For feature attention we use $g_i(\psi) =$
347 $1 + \beta h(\psi - \phi_i)$, where β the *feature attention gain*, and $h(\cdot)$ the *gain profile*. We follow the feature similarity
348 gain model (Treue and Martinez-Trujillo 1999), where a neuron's gain is enhanced if the attended feature
349 matches the neuron's preference and suppressed otherwise. A common choice for h is a cosine: $h(\psi - \phi_i) =$
350 $\cos(\psi - \phi_i)$.

351 Note that from the perspective of the model there is no fundamental difference between spatial and
352 feature attention. If we treat space as a variable that is being encoded by the population, any derivations
353 for feature attention also apply to spatial attention. However, because we consider only a local population
354 with overlapping receptive fields, spatial attention is a special case: the gain profile within the population
355 is constant and therefore spatial attention can be expressed in a simpler way using a single common
356 gain α . Thus, whenever we refer to spatial attention, this applies to a situation where all neurons in
357 the population that is being considered share the same preferred feature. Likewise, whenever we refer to

358 feature attention, this applies to any situation where the neurons in the population span a large range of
359 preferred features. We chose this (somewhat arbitrary) distinction, because it reflects the typical situation
360 in an experiment, where neurons with similar retinotopic locations are recorded, which typically span a
361 large range of preferred orientations or directions.

362 4.2 Effect of fluctuating gains on spike count statistics

Throughout this paper we assume that spatial and feature attention are independent processes and consider them in isolation. We further assume that the experimenter does not have access to the attentional state on individual trials, but can only control its average over many trials:

$$E[\alpha] = \mu \quad (27)$$

$$E[\beta] = \nu. \quad (28)$$

In addition the attentional state fluctuates from trial to trial with unknown variance

$$\text{Var}[\alpha] = \sigma^2 \quad (29)$$

$$\text{Var}[\beta] = \tau^2. \quad (30)$$

363 To compute means and (co-)variances of the observed spike counts we need only the means and variances
364 of α and β . The expected spike counts (Eqs. 3, 9) follow from the linearity of the expectation. Variances
365 and covariances can be computed by application of the Law of Total Variance (here for the case of spatial
366 attention, feature attention follows the same logic):

$$\text{Cov}[y_i, y_j] = E[\text{Cov}[y_i, y_j|\alpha]] + \text{Cov}[E[y_i|\alpha], E[y_j|\alpha]], \quad (31)$$

367 where the outer expectation (covariance) is taken over α and the inner covariance (expectation) over y_i
368 and y_j . Plugging the definitions of $\lambda_i = E[y_i|\alpha]$ and using the assumption of conditionally independent
369 Poisson spiking $\text{Cov}[y_i, y_j|\alpha] = \delta_{ij}\lambda_i$, we obtain the expressions for variances and covariances stated in
370 the main text (Eqs. 4–6, 10–13).

371 4.3 Effect of fluctuations in attended feature on spike count statistics

372 Calculating the means and covariances under fluctuations in the attended direction ψ follows the same
373 approach as above. However, since the gain profile $h_i(\psi)$ can be non-linear, we need a few additional
374 assumptions. We assume that ψ is distributed around some direction $\psi_0 = E[\psi]$ with variance $q^2 = \text{Var}[\psi]$.
375 For reasonably small q^2 we can approximate the gain profile by its first-order Taylor expansion

$$h_i(\psi) \approx h_i(\psi_0) + (\psi - \psi_0)h'_i(\psi_0), \quad (32)$$

376 where h'_i is the derivative with respect to ψ . Using this approximation we can write $E[h_i(\psi)] \approx h_i(\psi_0)$
377 and $\text{Var}[h_i(\psi)] \approx q^2 h'^2_i(\psi_0)$, which leads (again after applying the Law of Total Variance) to the results
378 in the main text (Eqs. 14–16).

379

4.4 Coding accuracy under fluctuations of spatial attention

380

Here we show that fluctuations in spatial attention have a negligible effect on the amount of information about the orientation of the stimulus. For simplicity we assume that neurons produce spikes conditionally independently given the stimulus orientation θ and the attentional gain g :

381

$$y_i|\theta, g \sim \text{Poisson}(\lambda_i) \quad \lambda_i = g f_i(\theta) \quad (33)$$

382

The attentional gain g is shared among all neurons and drawn from a Gamma distribution with shape μ^2/σ^2 and scale σ^2/μ , which implies $E[g] = \mu$ and $\text{Var}[g] = \sigma^2$. Assuming that the experimenter does not know the attentional gain, the distribution $P(\mathbf{y}|\theta)$ obtained by marginalizing over g is a multivariate negative binomial distribution:

$$P(\mathbf{y}|\theta) = \int P(g) \prod_i P(y_i|\theta, g) dg \quad (34)$$

$$= \int \frac{g^{\frac{\mu^2}{\sigma^2}-1} \exp(-\frac{g\mu}{\sigma^2})}{\Gamma(\frac{\mu^2}{\sigma^2}) (\frac{\sigma^2}{\mu})^{\frac{\mu^2}{\sigma^2}}} \prod_i \frac{(g f_i)^{y_i}}{y_i!} \exp(-g f_i) dg \quad (35)$$

$$= \frac{(\frac{\mu}{\sigma^2})^{\frac{\mu^2}{\sigma^2}}}{\Gamma(\frac{\mu^2}{\sigma^2})} \left(\prod_i \frac{f_i^{y_i}}{y_i!} \right) \int g^{\frac{\mu^2}{\sigma^2}-1+\sum y_i} \exp\left(-g\left(\frac{\mu}{\sigma^2} + \sum f_i\right)\right) dg \quad (36)$$

$$= \frac{\Gamma(\frac{\mu^2}{\sigma^2} + \sum y_i)}{\Gamma(\frac{\mu^2}{\sigma^2})} \left(\prod_i \frac{f_i^{y_i}}{y_i!} \right) \left(\frac{\mu}{\sigma^2} + \sum f_i \right)^{-\frac{\mu^2}{\sigma^2}} \left(\frac{1}{\frac{\mu}{\sigma^2} + \sum f_i} \right)^{\sum y_i} \quad (37)$$

For the Fisher information $J = E\left[\frac{d^2}{d\theta^2} \log P(\mathbf{y}|\theta)\right]$ we need the derivatives of the log-likelihood:

$$\frac{d}{d\theta} \log P(\mathbf{y}|\theta) = \left(\sum_i \frac{y_i f_i'}{f_i} \right) - \frac{(\frac{\mu^2}{\sigma^2} + \sum y_i) \sum f_i'}{\frac{\mu}{\sigma^2} + \sum f_i} \quad (38)$$

$$\frac{d^2}{d\theta^2} \log P(\mathbf{y}|\theta) = \left(\sum_i y_i \frac{f_i'' f_i - (f_i')^2}{f_i^2} \right) - \left(\frac{\mu^2}{\sigma^2} + \sum y_i \right) \frac{(\sum f_i'') (\frac{\mu}{\sigma^2} + \sum f_i) - (\sum f_i')^2}{(\frac{\mu}{\sigma^2} + \sum f_i)^2} \quad (39)$$

Plugging into the formula for Fisher information, re-ordering the summations over \mathbf{y} and i , and using the facts $\sum_{\mathbf{y}} P(\mathbf{y}|\theta) = 1$ and $\sum_{\mathbf{y}} P(\mathbf{y}|\theta) y_i = E[y_i] = \mu f_i$, we obtain

$$J = - \sum_{\mathbf{y}} P(\mathbf{y}|\theta) \left[\left(\sum_i y_i \frac{f_i'' f_i - (f_i')^2}{f_i^2} \right) - \left(\frac{\mu^2}{\sigma^2} + \sum y_i \right) \frac{(\sum f_i'') (\frac{\mu}{\sigma^2} + \sum f_i) - (\sum f_i')^2}{(\frac{\mu}{\sigma^2} + \sum f_i)^2} \right] \quad (40)$$

$$= \mu \sum_i \frac{(f_i')^2}{f_i} - \frac{\mu (\sum f_i')^2}{\frac{\mu}{\sigma^2} + \sum f_i} \quad (41)$$

383

The first term in the above equation is the Fisher information of an independent population of neurons and therefore $O(N)$, while the second term is $O(1)$: for homogeneous population of neurons, where $f_i(\theta) = f(\theta - \phi_i)$, it is zero; for heterogeneous populations it is $O(1)$, as we show in the next paragraph. Thus, fluctuations in spatial attention do not impair the coding accuracy of the population with respect to orientation.

384

385

386

387

388 To show that the second term above is $O(1)$ for heterogeneous populations, we assume that the
 389 neurons' tuning curves are independent random variables (see Ecker et al. 2011; Shamir and Sompolinsky
 390 2006). In this case the quantity of interest is the expected value with respect to different realizations of
 391 the heterogeneity:

$$E \left[\frac{\mu (\sum f'_i)^2}{\frac{\mu}{\sigma^2} + \sum f_i} \right] \approx \frac{\mu E[(\sum f'_i)^2]}{\frac{\mu}{\sigma^2} + E[\sum f_i]} = O(1). \quad (42)$$

392 Here the approximation holds because for large N the width of the distribution of $\sum f_i$ becomes narrower
 393 relative to its mean and therefore the expected value of the second term converges to the ratio of the
 394 expected values of numerator and denominator. The equality holds because $\sum f_i = O(N)$ and

$$E \left[(\sum f'_i)^2 \right] = \text{Var} \left[\sum f'_i \right] = \sum \text{Var}[f'_i] = O(N), \quad (43)$$

395 which holds because $E[\sum f'_i] = 0$.

396 4.5 Coding accuracy under fluctuations of feature attentional gain

397 Fluctuations in feature attention are more difficult to study analytically. Unfortunately, the Gamma-
 398 Poisson mixture model employed above does not generalize to the case where the gain is weighted dif-
 399 ferently for each neuron (i.e. the gain profile h_i), or at least we are not aware of a model that has a
 400 closed-form expression for the marginal probability mass function when the gain is unknown. Therefore,
 401 we here approximate the population activity by a multivariate Gaussian distribution with matching mean
 402 and covariance matrix (Eq. 9–13) and focus on linear readout. Under this approximation, the (linear)
 403 Fisher Information is given by

$$J = (\mathbf{f}')^T C^{-1} \mathbf{f}'. \quad (44)$$

404 The inverse of the covariance matrix is obtained by applying a rank-one update:

$$C^{-1} = F^{-1} - \frac{F^{-1} \mathbf{u} \mathbf{u}^T F^{-1}}{\tau^{-2} + \mathbf{u}^T F^{-1} \mathbf{u}}, \quad (45)$$

where $F_{ii} = (1 + \nu h_i(\psi)) f_i(\theta)$ and $u_i = h_i(\psi) f_i(\theta)$ as above. Plugging in and simplifying we obtain

$$J = J_0 - \frac{\left(\sum \frac{h_i f'_i}{1 + \nu h_i} \right)^2}{\tau^{-2} + \sum \frac{h_i f_i}{1 + \nu h_i}} \quad (46)$$

$$= J_0 - O(1). \quad (47)$$

405 As above for spatial attention, the $O(1)$ correction term is exactly zero for homogeneous populations and
 406 the derivation for heterogeneous populations follows the same line of argument as above.

407 4.6 Coding accuracy under fluctuations of attended feature

As described in the main text, fluctuations of the attended feature create differential correlations, i.e.
 response variability that is identical to variability induced by changes in the stimulus. Here we derive

the results using the Generalized Linear Model (Eq. 25) formulation:

$$\log \lambda_i = \beta \cos(\psi - \phi_i) + \kappa \cos(\theta - \phi_i) \quad (48)$$

$$= (\mathbf{b} + \kappa \mathbf{x})^T \mathbf{k}_i \quad (49)$$

$$\equiv \hat{\mathbf{x}}^T \mathbf{k}_i, \quad (50)$$

408 where $\mathbf{b} = \beta[\cos \psi, \sin \psi]^T$, $\mathbf{x} = [\cos \theta, \sin \theta]^T$, and $\mathbf{k}_i = [\cos \phi_i, \sin \phi_i]^T$. Since $\hat{\mathbf{x}}$ is independent of the
 409 neurons, it is obvious that attention has exactly the same effect as a change in the stimulus. Assuming
 410 $E[\psi] = \theta$, $\text{Var}[\psi]$ is small, and (without loss of generality) $\theta = 0$, we have

$$\mathbf{b} \approx \beta \begin{bmatrix} 1 \\ \psi \end{bmatrix}, \quad \mathbf{x} \approx \begin{bmatrix} 1 \\ 0 \end{bmatrix}. \quad (51)$$

411 Moreover, we can write the attention-perturbed stimulus $\hat{\theta}$ as

$$\hat{\theta} \approx \frac{\hat{x}_2}{\hat{x}_1} = \frac{\psi}{\kappa + \beta}. \quad (52)$$

412 For large N the Poisson noise averages out and therefore the resulting Fisher information is simply the
 413 inverse of the variance of the (attention-perturbed) stimulus:

$$J \rightarrow \frac{1}{\text{Var}[\hat{\theta}]} = \frac{(\kappa/\beta + 1)^2}{\text{Var}[\psi]}. \quad (53)$$

References

- 414
- 415 Abbott, L. F. and Peter Dayan (1999). “The Effect of Correlated Variability on the Accuracy of a
416 Population Code”. In: *Neural Computation* 11.1. 00475, pp. 91–101.
- 417 Bair, Wyeth, Ehud Zohary, and William T. Newsome (2001). “Correlated Firing in Macaque Visual
418 Area MT: Time Scales and Relationship to Behavior”. In: *The Journal of Neuroscience* 21.5. 00299,
419 pp. 1676–1697.
- 420 Britten, Kenneth H. et al. (1993). “Responses of neurons in macaque MT to stochastic motion signals”.
421 In: *Visual Neuroscience* 10.06, pp. 1157–1169.
- 422 Buesing, Lars, Jakob H. Macke, and Maneesh Sahani (2012). “Learning stable, regularised latent models
423 of neural population dynamics”. In: *Network: Computation in Neural Systems* 23.1-2. 00003, pp. 24–
424 47.
- 425 Cohen, Marlene R and John H R Maunsell (2009). “Attention improves performance primarily by reducing
426 interneuronal correlations”. In: *Nature Neuroscience* 12.12, pp. 1594–1600. ISSN: 1097-6256.
- 427 Cohen, Marlene R. and John H. R. Maunsell (2010). “A Neuronal Population Measure of Attention Pre-
428 dictors Behavioral Performance on Individual Trials”. In: *The Journal of Neuroscience* 30.45, pp. 15241–
429 15253.
- 430 Cohen, Marlene R. and John H.R. Maunsell (2011). “Using Neuronal Populations to Study the Mecha-
431 nisms Underlying Spatial and Feature Attention”. In: *Neuron* 70.6, pp. 1192–1204. ISSN: 0896-6273.
- 432 Collins, Michael, Sanjoy Dasgupta, and Robert E. Schapire (2001). “A generalization of principal com-
433 ponents analysis to the exponential family”. In: *Advances in neural information processing systems*.
434 00259, pp. 617–624.
- 435 Dean, A. F. (1981). “The variability of discharge of simple cells in the cat striate cortex”. en. In: *Experi-
436 mental Brain Research* 44.4, pp. 437–440. ISSN: 0014-4819, 1432-1106.
- 437 Ecker, Alexander S. and Andreas S. Tolias (2014). “Is there signal in the noise?” en. In: *Nature Neuro-
438 science* 17.6. 00000, pp. 750–751. ISSN: 1097-6256.
- 439 Ecker, Alexander S. et al. (2010). “Decorrelated Neuronal Firing in Cortical Microcircuits”. In: *Science*
440 327.5965, pp. 584–587.
- 441 Ecker, Alexander S. et al. (2011). “The Effect of Noise Correlations in Populations of Diversely Tuned
442 Neurons”. In: *The Journal of Neuroscience* 31.40. 00032, pp. 14272 –14283.
- 443 Ecker, Alexander S. et al. (2014). “State dependence of noise correlations in macaque primary visual
444 cortex”. In: *Neuron* 82.1, pp. 235–248.
- 445 Goris, Robbe L. T., J. Anthony Movshon, and Eero P. Simoncelli (2014). “Partitioning neuronal variabil-
446 ity”. en. In: *Nature Neuroscience* 17.6. 00000, pp. 858–865. ISSN: 1097-6256.
- 447 Herrero, Jose L. et al. (2013). “Attention-Induced Variance and Noise Correlation Reduction in Macaque
448 V1 Is Mediated by NMDA Receptors”. In: *Neuron* 78.4, pp. 729–739. ISSN: 0896-6273.
- 449 Josić, Krešimir et al. (2009). “Stimulus-Dependent Correlations and Population Codes”. In: *Neural Com-
450 putation* 21.10. 00038, pp. 2774–2804. ISSN: 0899-7667.
- 451 Kay, Steven M. (1993). *Fundamentals of Statistical Signal Processing, Volume I: Estimation Theory*.
452 1st ed. Prentice Hall. ISBN: 0133457117.
- 453 Macke, Jakob H. et al. (2011). “Empirical models of spiking in neural populations”. In: *Advances in neural
454 information processing systems* 24, p. 13501358.
- 455 Maunsell, John H.R. and Stefan Treue (2006). “Feature-based attention in visual cortex”. In: *Trends in
456 Neurosciences* 29.6, pp. 317–322. ISSN: 0166-2236.

- 457 Mitchell, Jude F., Kristy A. Sundberg, and John H. Reynolds (2009). “Spatial Attention Decorrelates
458 Intrinsic Activity Fluctuations in Macaque Area V4”. In: *Neuron* 63.6, pp. 879–888. ISSN: 0896-6273.
- 459 Mohamed, Shakir, Zoubin Ghahramani, and Katherine A. Heller (2009). “Bayesian exponential family
460 PCA”. In: *Advances in Neural Information Processing Systems*. 00037, pp. 1089–1096.
- 461 Moreno-Bote, Rubén et al. (2014). “Information-limiting correlations”. en. In: *Nature Neuroscience* ad-
462 vance online publication. 00000. ISSN: 1097-6256.
- 463 Nienborg, Hendrikje and Bruce G. Cumming (2009). “Decision-related activity in sensory neurons reflects
464 more than a neuron’s causal effect”. In: *Nature* 459.7243, pp. 89–92. ISSN: 0028-0836.
- 465 Reynolds, John H. and Leonardo Chelazzi (2004). “Attentional Modulation of Visual Processing”. In:
466 *Annual Review of Neuroscience* 27.1, pp. 611–647. ISSN: 0147-006X.
- 467 Rocha, Jaime de la et al. (2007). “Correlation between neural spike trains increases with firing rate”. In:
468 *Nature* 448.7155, pp. 802–806. ISSN: 0028-0836.
- 469 Shadlen, Michael N. and William T. Newsome (1998). “The Variable Discharge of Cortical Neurons: Im-
470 plications for Connectivity, Computation, and Information Coding”. In: *The Journal of Neuroscience*
471 18.10. 01334, pp. 3870–3896.
- 472 Shamir, Maoz and Haim Sompolinsky (2006). “Implications of Neuronal Diversity on Population Coding”.
473 In: *Neural Computation* 18.8, pp. 1951–1986.
- 474 Smith, Matthew A. and Adam Kohn (2008). “Spatial and Temporal Scales of Neuronal Correlation in
475 Primary Visual Cortex”. In: *J. Neurosci.* 28.48, pp. 12591–12603.
- 476 Smith, Matthew A. and Marc A. Sommer (2013). “Spatial and Temporal Scales of Neuronal Correlation
477 in Visual Area V4”. en. In: *The Journal of Neuroscience* 33.12. 00007, pp. 5422–5432. ISSN: 0270-6474,
478 1529-2401.
- 479 Sompolinsky, Haim et al. (2001). “Population coding in neuronal systems with correlated noise”. In:
480 *Physical Review E* 64.5, p. 051904.
- 481 Tolhurst, David J., J. A. Movshon, and A. F. Dean (1983). “The statistical reliability of signals in single
482 neurons in cat and monkey visual cortex”. In: *Vision research* 23.8. 00658, pp. 775–785.
- 483 Treue, Stefan and Julio C. Martinez-Trujillo (1999). “Feature-based attention influences motion processing
484 gain in macaque visual cortex”. In: *Nature* 399.6736, pp. 575–579. ISSN: 0028-0836.
- 485 Wilke, Stefan D. and Christian W. Eurich (2002). “On the functional role of noise correlations in the
486 nervous system”. In: *Neurocomputing* 44-46, pp. 1023–1028. ISSN: 0925-2312.
- 487 Zohary, Ehud, Michael N. Shadlen, and William T. Newsome (1994). “Correlated neuronal discharge rate
488 and its implications for psychophysical performance”. In: *Nature* 370.6485. 00736, pp. 140–143.