

1 Attentional modulation of neuronal 2 variability in circuit models of cortex

3 Tatjana Kanashiro^{1,2,3}, Gabriel Koch Ocker^{2,3,4}, Marlene R. Cohen^{3,5}, Brent
4 Doiron^{2,3}

*For correspondence:
bdoiron@pitt.edu (BD)

5 ¹Program for Neural Computation, Carnegie Mellon University and University of Pittsburgh,
6 Pittsburgh, PA, USA; ²Department of Mathematics, University of Pittsburgh, Pittsburgh, PA,
7 USA; ³Center for the Neural Basis of Cognition, Pittsburgh, PA, USA; ⁴Allen Institute for
8 Brain Science, Seattle, WA, USA; ⁵Department of Neuroscience, University of Pittsburgh,
9 Pittsburgh, PA, USA

10

11 **Abstract** The circuit mechanisms behind shared neural variability (noise correlation) and its
12 dependence on neural state are poorly understood. Visual attention is well-suited to constrain cortical
13 models of response variability because attention both increases firing rates and their stimulus
14 sensitivity, as well as decreases noise correlations. We provide a novel analysis of population
15 recordings in rhesus primate visual area V4 showing that a single biophysical mechanism may underlie
16 these diverse neural correlates of attention. We explore model cortical networks where top-down
17 mediated increases in excitability, distributed across excitatory and inhibitory targets, capture the key
18 neuronal correlates of attention. Our models predict that top-down signals primarily affect inhibitory
19 neurons, whereas excitatory neurons are more sensitive to stimulus specific bottom-up inputs.
20 Accounting for trial variability in models of state dependent modulation of neuronal activity is a
21 critical step in building a mechanistic theory of neuronal cognition.

22

23 Introduction

24 The behavioral state of the brain exerts a powerful influence on the cortical responses. For example,
25 electrophysiological recordings from both rodents and primates show that the level of wakefulness
26 (*Steriade et al., 1993*), active sensory exploration (*Crochet et al., 2011*), and attentional focus
27 (*Treue, 2001; Reynolds and Chelazzi, 2004; Gilbert and Sigman, 2007; Moore and Zirnsak,
28 2017*) all modulate synaptic and spiking activity. Despite the diversity of behavioral contexts, in all of
29 these cases an overall elevation and desynchronization of cortical activity accompanies heightened states
30 of processing (*Harris and Thiele, 2011*). Exploration of the neuronal mechanisms that underly
31 such state changes has primarily centered around how various neuromodulators shift the cellular and
32 synaptic properties of cortical circuits (*Hasselmo, 1995; Lee and Dan, 2012; Noudoost and
33 Moore, 2011; Moore and Zirnsak, 2017*). However, a coherent theory linking the modulation of
34 cortical circuits to an active desynchronization of population activity is lacking. In this study we
35 provide a circuit-based theory for the known attention-guided modulations of neuronal activity in the
36 visual cortex of primates performing a stimulus change detection task.

37 The investigation of the neuronal correlates of attention has a rich history. Attention increases
38 the firing rates of neurons engaged in feature- and spatial-based processing tasks (*McAdams and
39 Maunsell, 2000; Reynolds et al., 1999*). Attentional modulation of the stimulus-response sensi-
40 tivity (gain) of firing rates is more complicated, often depending on stimulus specifics such as the size
41 and contrast of a visual image (*Williford and Maunsell, 2006; Reynolds and Heeger, 2009;
42 Sanayei et al., 2015*). In recent years there has been increased focus on how brain states affect

43 trial-to-trial spiking variability (*Crochet et al., 2011; Lin et al., 2015; Doiron et al., 2016;*
 44 *Stringer et al., 2016*). In particular, attention decreases the shared variability (noise correlations) of
 45 the firing rates from pairs of neurons (*Cohen and Maunsell, 2009; Mitchell et al., 2009; Cohen*
46 and Maunsell, 2011; Herrero et al., 2013; Ruff and Cohen, 2014; Engel et al., 2016). The
 47 combination of a reduction in noise correlations and an increase in response gain has potentially
 48 important functional consequences through an improved population code (*Cohen and Maunsell,*
49 2009; Rabinowitz et al., 2015). In total, there is an emerging picture of the impact of attention
 50 on the trial-averaged and trial-variable spiking dynamics of cortical populations.

51 Phenomenological models of attentional modulation have been popular (*Reynolds and Heeger,*
52 Navalpakkam and Itti, 2005; Gilbert and Sigman, 2007; Ecker et al., 2016); however,
 53 such analyses cannot provide insight into the circuit mechanics of attentional modulation. Biophysical
 54 models of attention circuits are difficult to constrain, due in large part to the diversity of mechanisms
 55 which control the firing rate and response gain of neurons (*Silver, 2010; Sutherland et al., 2009*).
 56 Nonetheless, several circuit models for attentional modulation have been proposed (*Ardid et al.,*
57 Deco and Thiele, 2011; Buia and Tiesinga, 2008), but analysis has been mostly confined
 58 to trial-averaged responses. Taking inspiration from these studies, mechanistic models of attentional
 59 modulation can be broadly grouped along two hypotheses. First, the circuit mechanisms that control
 60 trial-averaged responses (i.e firing rates and response gain) may be distinct from those that modulate
 61 neuronal variability (i.e noise correlations). This hypothesis has support from experiments in primate
 62 V1 showing that N-methyl-D-aspartate receptors have no impact on top-down attentional modulation
 63 of firing rates, yet have a strong influence of attentional control of noise correlations (*Herrero et al.,*
64 2013). A second hypothesis is that the modulations of firing rates and noise correlations are reflections
 65 of a single biophysical mechanism. Support for this comes from pairs of V4 neurons that each show
 66 strong attentional modulation of firing rates, also show a strong attention mediated reductions in
 67 noise correlation (*Cohen and Maunsell, 2011*). In this study we provide novel analysis of the
 68 covariability of V4 population activity engaged in an attention-guided detection task (*Cohen and*
69 Maunsell, 2009) that is consistent with the second hypothesis. Specifically, the modulation of spike
 70 count covariance between unattended and attended states has the same dimensionality as the firing
 71 rate modulation.

72 We use the results from our dimensionality analysis to show that an excitatory-inhibitory recurrent
 73 circuit model subject to global fluctuations is sufficient to capture both the increase in firing rate
 74 and response gain as well as population-wide decrease of noise correlations. Our model makes two
 75 predictions regarding neuronal modulation: 1) that attentional modulation favors inhibitory neurons,
 76 and 2) that stimulus drive favors excitatory neurons. Finally, we show that our model predicts
 77 increased informational content in the excitatory population, which would result in improved readout
 78 by potential downstream targets. In total, our study provides a simple, parsimonious, and biologically
 79 motivated model of attentional modulation in cortical networks.

80 Results

81 Attention decreases noise correlations primarily by decreasing covariance

82 Two rhesus monkeys (*Macaca mulatta*) with microelectrode arrays implanted bilaterally in V4 were
 83 trained in an orientation change detection task (Fig. 1a; see Methods: Data preparation). A display
 84 with oriented Gabor gratings on the left and right flashed on and off. The monkey was cued to attend
 85 to either the left or right grating before each block of trials, while keeping fixation on a point between
 86 the two gratings. After a random number of presentations, one of the gratings changed orientation.
 87 The monkey then had to saccade to that side to obtain a reward. The behavioral task and data
 88 collection have been previously reported (*Cohen and Maunsell, 2009*).

89 A neuron is considered to be in an “attended state” when the attended stimulus is in the hemifield
 90 containing that neuron’s receptive field (contralateral hemifield), and in an “unattended state” when
 91 it is in the other (ipsilateral) hemifield. The trial-averaged firing rates from both attended and
 92 unattended neurons displayed a brief transient rise (~100 ms after stimulus onset), and eventually

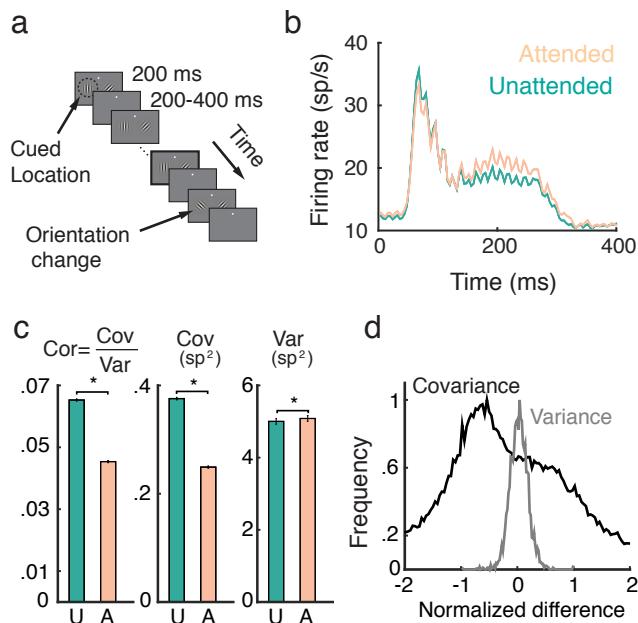


Figure 1. Attention increases firing rates and decreases trial-to-trial covariability of population responses. **a**, Overview of orientation-change detection task. See *Cohen and Maunsell (2009)* for a full description. **b**, Firing rates of neurons in the unattended (turquoise) and attended (orange) states, averaged over 3,170 units. The slight oscillation in the firing rate was due to the monitor refresh rate. **c**, Attention significantly decreased the spike count correlation and covariance and slightly increased variance. Error bars provide the SEM. **d**, Histograms of changes in covariance for each unit pair (black) and variance for each unit (gray). In each case we consider the relative change $[X^A - X^U]/\max(X^A, X^U)$, where X is either $\text{Cov}(n_i, n_j)$ or $\text{Var}(n_i)$. Data was collected from two monkeys over 21 and 16 recording sessions respectively. Signals were analyzed over a 200 ms interval, starting 60 ms after stimulus onset.

settled to an elevated sustained rate before the trial concluded (Fig. 1b). During the sustained period the mean firing rate of attended neurons (22.0 sp/s) was greater than that of unattended neurons (20.6 sp/s) (t test, $P < 10^{-5}$).

A major finding of *Cohen and Maunsell (2009)* was that the pairwise trial-to-trial noise correlations of the neuronal responses decreased with attention (Fig. 1c, left, mean unattended .065, mean attended .045, t test, $P < 10^{-5}$). The noise correlation between neurons i and j is a normalized measure, $\rho_{ij} = \text{Cov}(n_i, n_j)/\sqrt{\text{Var}(n_i)\text{Var}(n_j)}$, where Cov and Var denote spike count covariance and variance respectively. Both spike count variance and covariance significantly change with attention ($\langle \text{Var}^U \rangle_{\text{trials}} = 5.02 \text{ spikes}^2$, $\langle \text{Var}^A \rangle_{\text{trials}} = 5.10 \text{ spikes}^2$, t test, $P < 10^{-3}$; $\langle \text{Cov}^U \rangle_{\text{trials}} = 0.379 \text{ spikes}^2$, $\langle \text{Cov}^A \rangle_{\text{trials}} = 0.252 \text{ spikes}^2$, t test, $P < 10^{-5}$), but the decrease in covariance (34.0%) is much more pronounced than the increase in variance (1.61%; Fig. 1c, middle and right). We therefore conclude that the attention mediated decrease in noise correlation is primarily due to decreased covariance.

To further validate this observation, we consider the distributions of pairwise changes in covariance (black) and variance (gray) with attention over the entire data set (Fig. 1d). Covariance and variance are normalized by their respective maximal unattended or attended values (see Methods: Comparing change in covariance to change in variance). The change in covariance with attention is concentrated below zero with a large spread, whereas the change in variance is centered on zero with a narrower spread. Taken together these results suggest that to understand the mechanism by which noise correlations decrease it is necessary and sufficient to understand how spike count covariance decreases with attention.

114 Attention is a low-rank modulation of noise covariance

115 A reasonable simplification of V4 neurons is that they receive a bottom-up stimulus alongside an
116 attention-mediated top-down modulatory input. However, to properly model top-down attention we
117 need to first understand the dimension of attentional modulation on the V4 circuit as a whole. Let
118 $A_\phi : \phi^U \mapsto \phi^A$ denote the attentional modulation of measure ϕ from its value in the unattended
119 state, ϕ^U , to its value in the attended state, ϕ^A . For example, the firing rate modulation A_r can be
120 written as $\mathbf{r}^A = A_r \circ \mathbf{r}^U$, where \mathbf{r}^A is an $N \times 1$ vector of neural firing rates in the attended state,
121 \mathbf{r}^U denotes the firing rate vector in the unattended state, A_r is a vector the same size as \mathbf{r} , and \circ
122 denotes elementwise multiplication. In this case, the entries a_i of A_r are the ratios of the firing rates:
123 $a_i = r_i^A / r_i^U$ (Fig. 2a).

124 A less trivial aspect of attentional modulation is the modulation of covariance matrices:

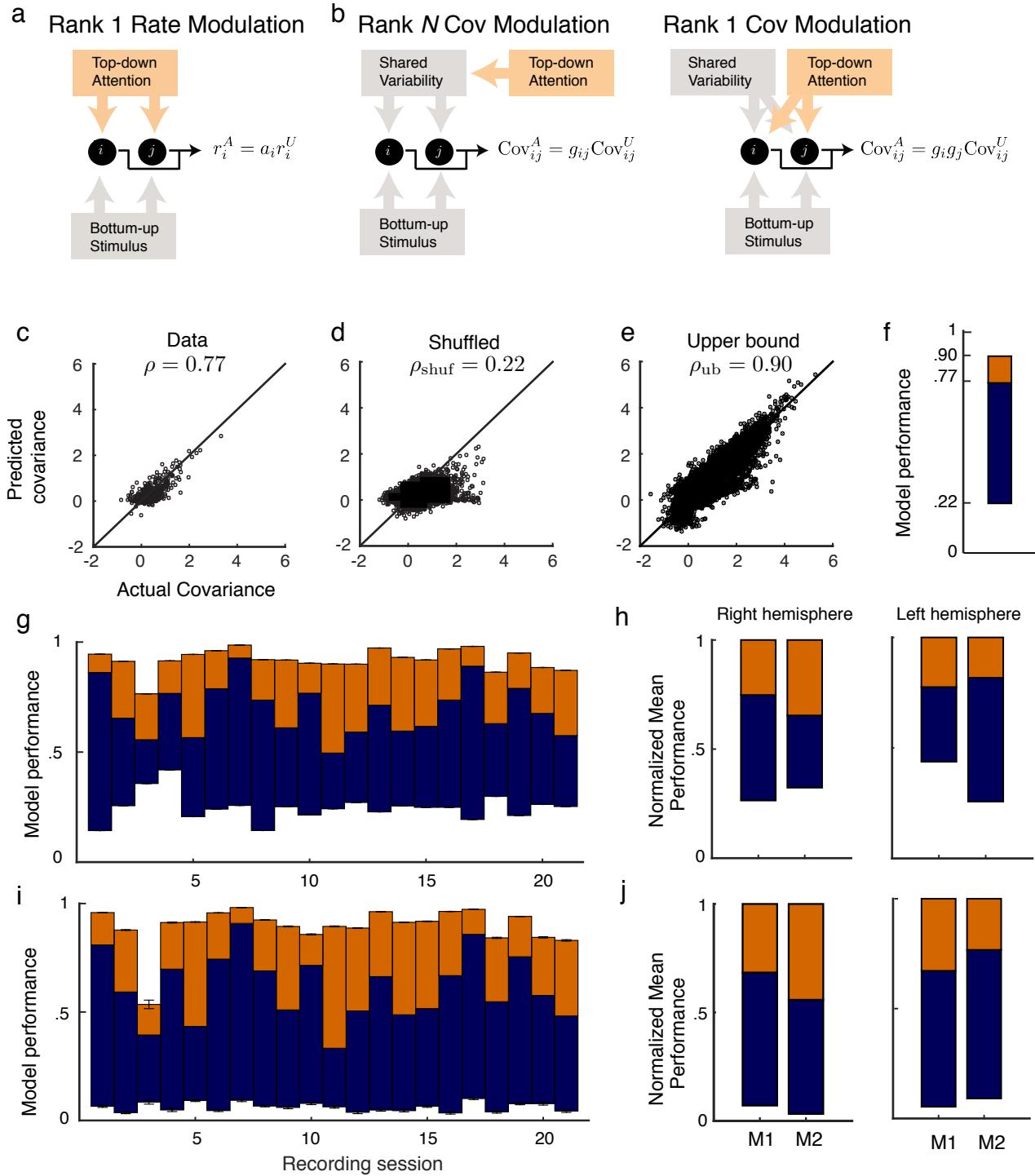
$$\mathbf{C}^A = A_C \circ \mathbf{C}^U. \quad (1)$$

125 Here \mathbf{C}^A is the attended spike count covariance matrix, \mathbf{C}^U the unattended spike count covariance
126 matrix, and A_C is a matrix the same size as \mathbf{C}^U , consisting of entries g_{ij} , which we will call *covariance*
127 *gains*. Unlike firing rates, the transformation matrix A_C can be of varying rank. On the one hand
128 A_C could be constructed from the ratios of the individual elements: $g_{ij} = c_{ij}^A / c_{ij}^U$, with each pair of
129 neurons (i, j) receiving an individualized attentional modulation g_{ij} of their shared variability (Fig.
130 2b, left). Under this modulation A_C is a rank N matrix. A rank N A_C will always perfectly (and
131 trivially) capture the matrix mapping in Eq. (1). However, it is difficult to conceive of a top-down
132 circuit mechanism that would allow attention to modulate each pair individually. On the other hand,
133 g_{ij} could depend not on the specific pair (i, j) , but on the individual neurons of the pairing: $g_{ij} = g_i g_j$
134 (Fig. 2b, right). In this case, only N values are needed to characterize A_C : $A_C = \mathbf{g} \mathbf{g}^T$, where \mathbf{g} is
135 a $N \times 1$ column vector, meaning A_C has rank of 1. This is a more parsimonious and biophysically
136 plausible scenario for attentional modulation, since in this case the covariance gain g_{ij} of neurons i
137 and j is simply emergent from the attentional modulation of the individual neurons. To test whether
138 A_C is low rank we analyzed the V4 population recordings during the visual attention task (Fig. 1),
139 specifically measuring A_C under the assumption that A_C is rank 1:

$$\mathbf{C}^A = \mathbf{g} \mathbf{g}^T \circ \mathbf{C}^U. \quad (2)$$

140 Equation (2) is a system of $N(N - 1)/2$ equations of the form $c_{ij}^A = g_i g_j c_{ij}^U$ in N unknowns
141 $\mathbf{g} = [g_1, \dots, g_N]^T$ (we only consider $i \neq j$ to exclude variance modulation from our analysis). For $N > 3$
142 this is an overdetermined system, and we solve for \mathbf{g} using a nonlinear equation solver. Let $\hat{\mathbf{g}}$ be the
143 optimal solution obtained by the solver (measured as a minimization of the L^2 -norm of the error; see
144 Methods:). Then $\hat{\mathbf{C}}^A := \hat{\mathbf{g}} \hat{\mathbf{g}}^T \circ \mathbf{C}^U$ provides an approximation to the attended covariance matrix.
145 In an example data set from a single recording session with $N = 39$ units, the correlation coefficient
146 ρ of the actual attended covariance values from \mathbf{C}^A versus the approximated attended covariance
147 values from $\hat{\mathbf{C}}^A$ was 0.77 (Fig. 2c). A shuffled \mathbf{C}^A matrix provides a reasonable null model, and the
148 example data set produces the lower bound correlation $\rho_{\text{shuf}} = 0.22$ (Fig. 2d; see Methods: Shuffled
149 covariance matrices). Finally, a Poisson model that perfectly decomposes as Eq. (2), yet sampled
150 with the same number of trials as in the experiment, gives an upper bound for the rank one structure,
151 the example data yields $\rho_{\text{ub}} = 0.90$ (Fig. 2e; see Methods: Upper bound covariance matrices). In
152 total, the combination of ρ , ρ_{shuf} , and ρ_{ub} (Fig. 2f) suggests that the rank one model of attention
153 modulation of covariance A_C is well justified.

154 We applied this analysis to 21 recording sessions from the right hemisphere of one monkey (Fig.
155 2g). For most of the recording sessions ρ is closer to ρ_{ub} than ρ_{shuf} . The averaged performance of
156 all sessions for both hemispheres of two monkeys generally agreed with this trend (Fig. 2h) We
157 normalized ρ and ρ_{shuf} by ρ_{ub} for each session to better compare different sessions that were subject
158 to day-to-day variations outside of the experimenter's control, such as the task performance or the
159 internal state of the monkey. To further validate our model we show the distribution of g_i s computed
160 from the entire data set (Fig. 3a). The majority of g_i values are less than one, consistent with

**Figure 2.** Caption is on next page.

161 $\langle \text{Cov}^A \rangle_{\text{trials}} < \langle \text{Cov}^U \rangle_{\text{trials}}$ (Fig. 1c). Further, there was little relation between the attentional
162 modulation of firing rates, measured by r_i^A / r_i^U , and the attentional modulation of covariance through
163 g_i (Fig. 3b). This indicates that the circuit modulation of firing rates and covariance are not be
164 trivially related to one another (Doiron *et al.*, 2016).

Figure 2. Rank one structure of attentional modulation of spike count covariance. **a**, Attentional modulation of firing rate. Firing rates of neurons i and j (black circles are modulated by bottom-up stimulus and top-down attention. **b**, Two possible models of attentional modulation of covariance. Left: High-rank covariance modulation, in which attention modulates the shared variability of each pair of neurons. Right: Low-rank covariance modulation, in which attention modulates each neuron individually rather than in a pairwise manner. **c-e**, The measured covariance values plotted against those predicted by the rank-1 model for data collected in one recording session, for **c**, the actual data ($\rho = 0.77$), **d**, shuffled data ($\rho_{\text{shuf}} = 0.22$, 100 shuffles), and **e**, artificial upper-bound data ($\rho_{\text{ub}} = 0.90$, 10 realizations of the upper bound model). **f**, Synthesis of **c-e** in a bar plot. The orange area represents the loss of model performance compared to the upper bound model, and the blue area represents the increase in model performance compared to model applied to shuffled data. **g**, Rank-1 model performance reported for 21 recording sessions from one monkey. Each bar represents one recording session. Recordings from a mean of $N = 53.5$ units in the right-hemisphere were analyzed, with maximum and minimum N of 80 and 35, respectively. Error bars denote standard error of the mean. **h**, Mean normalized performance (relative to ρ_{ub}) for both hemispheres of two monkeys (M1 and M2). **i**, Analysis as in **g**, using leave-one-out cross-validation to test the predictive power of the model. **j**, Mean normalized performance of the cross-validated data.

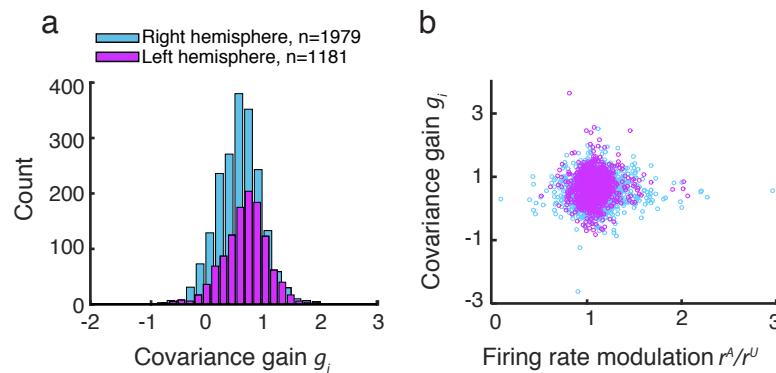


Figure 3. Covariance gain shows the attenuation of population-wide fluctuations with attention. **a**, Distribution of covariance gains g_i computed from the entire data set. **b**, The relation between covariance g_i and the attention mediated modulation of firing rates r_i^A/r_i^U . The correlation coefficients between the data sets were 0.036 and 0.051 for the right and left hemispheres, respectively.

165 We additionally tested the validity of our model in Eq. (2) with a leave-one-out cross-validation
 166 analysis (see Methods: Leave-one-out cross-validation). We accurately predicted an omitted covariance
 167 C_{ij}^A (Fig. 2i and j), consistent with our original analysis (Fig. 2g and h). The individual session-by-
 168 session performance values for both the standard and leave-one-out setups are provided (Appendix:
 169 Model performance for all monkeys and hemispheres).

170 Finally, we investigated to what extent the actual value of the covariance gain g_i of neuron i
 171 depends on the population of neurons in which it was computed. We solved the system of equations
 172 $C_{ij}^A = g_i g_j C_{ij}^U$ using covariance matrices computed from recordings from distinct sets of neurons,
 173 overlapping only by neuron i . This gives two estimates of g_i , that nevertheless agreed largely with
 174 one another (Appendix: Low-dimensional modulation is intrinsic to neurons). This supported the
 175 hypothesis that covariance gain g_i is an intrinsic property of neuron i .

176 The standard and cross-validation tests verify that the low-rank model of attentional modulation
 177 defined in Eq. (2) explains between 66 and 82% (standard), or 56 and 77% (cross-validation) of the
 178 data. Taking this to be a positive result, we conclude that the covariance gain modulation depends
 179 largely on the modulation of individual neurons.

180 Network requirements for attentional modulation

181 Having described attentional modulation statistically our next goal is to develop a circuit model
 182 to understand the process mechanistically. Consider a network of N coupled neurons, and let the
 183 spike count from neuron i on a given trial be y_i . The network output has the covariance matrix \mathbf{C}

184 with elements $c_{ij} = \text{Cov}(y_i, y_j)$. In this section we identify the minimal circuit elements so that the
185 attentional mapping $A_C : \mathbf{C}^U \mapsto \mathbf{C}^A$ satisfies the following two conditions (on average):

186 **C1:** $c_{ij}^A = g_i g_j c_{ij}^U$; attentional modulation of covariance is rank one (Fig. 2).

187 **C2:** $g_i < 1$; spike count covariance decreases with attention (Fig. 1).

188 What follows is only a sketch of our derivation (a complete treatment is given in Appendix: Network
189 requirements for attentional modulation).

190 If inputs are weak then y_i can be described by a linear perturbation about a background
191 state (*Ginzburg and Sompolinsky, 1994; Doiron et al., 2004; Trousdale et al., 2012*):

$$y_i = y_{iB} + L_i \left(\sum_{k=1}^N J_{ik} y_k + \xi_i \right). \quad (3)$$

192 Here y_{iB} is the background activity of neuron i , J_{ik} is the coupling strength from neuron k to i , and
193 L_i is the input-to-output gain of neuron i . In addition to internal coupling we assume a source of
194 external fluctuations ξ_i to neuron i . Here y_i , y_{iB} , and ξ_i are random variables that vary across trials.
195 The trial-averaged firing rate of neuron i is $r_i = \langle y_i \rangle / T$ (where $\langle \cdot \rangle$ denotes averaging over trials of
196 length T). The background state has variability $b_i = \text{Var}(y_{iB})$ which we assume to be independent
197 across neurons, meaning the background network covariance is $\mathbf{B} = \text{diag}(b_i)$. Finally, the external
198 fluctuations have covariance matrix \mathbf{X} with element $x_{ij} = \text{Cov}(\xi_i, \xi_j)$.

199 Motivated by our analysis of population recordings (Fig. 2) we study attentional modulations
200 that target individual neurons. This amounts to considering only $A_r : r_i^U \mapsto r_i^A$ and $A_L : L_i^U \mapsto L_i^A$.
201 Additionally, we assume that any model of attentional modulation must result in $r_i^A > r_i^U$ (Fig. 1b).
202 A widespread property of both cortical pyramidal cells and interneurons is that an increase of firing
203 rate r_i causes an increase of input-output gain L (*Cardin et al., 2007*), thus we will also require
204 $L^A > L^U$.

205 Spiking covariability in recurrent networks can be due to internal interactions (through J_{ik}) or
206 external fluctuations (through ξ_i), or both (*Ocker et al., 2017*). Networks with unstructured
207 connectivity have internally generated covariability that vanishes as N grows. This is true if the
208 connectivity is sparse (*van Vreeswijk and Sompolinsky, 1998*), or dense having weak synapses
209 where $J_{ik} \sim 1/N$ (*Trousdale et al., 2012*) or strong synapses where $J_{ik} \sim 1/\sqrt{N}$ combined with
210 a balance between excitation and inhibition (*Renart et al., 2010; Rosenbaum et al., 2017*). In
211 these cases spiking covariability requires external fluctuations to be applied and subsequently filtered
212 by the network. We follow this second scenario and choose \mathbf{X} so as to provide external covariability to
213 our network.

214 Recent analysis of cortical population recordings show that the shared spiking variability across
215 the population can be well approximated by a rank one model of covariability (*Kelly et al., 2010;*
216 *Ecker et al., 2014; Lin et al., 2015; Ecker et al., 2016; Rabinowitz et al., 2015; Whiteway*
217 *and Butts, 2017*) (we remark that *Rabinowitz et al. (2015)* analyzed the same data set that we
218 have in Figs. 1 and 2). Thus motivated we take the external fluctuations \mathbf{X} to be rank one with
219 $x_{ij} = x_i x_j$, reflecting a single source of global external variability ξ with unit variance (neuron i
220 receives $\xi_i = x_i \xi$). Combining this assumption with the linear ansatz in Eq. (3) and taking N large
221 yields:

$$\mathbf{C} \approx ((\mathbf{I} - \mathbf{K})^{-1} \mathbf{L} \mathbf{x}) ((\mathbf{I} - \mathbf{K})^{-1} \mathbf{L} \mathbf{x})^T = \mathbf{c} \mathbf{c}^T, \quad (4)$$

222 where matrix \mathbf{K} has element $K_{ij} = L_i J_{ij}$ and $\mathbf{L} = \text{diag}(L_i)$. We have also defined the vectors
223 $\mathbf{x} = [x_1, \dots, x_N]^T$ and $\mathbf{c} = [c_1, \dots, c_N]^T$ with $c_i = ((\mathbf{I} - \mathbf{K})^{-1} \mathbf{L} \mathbf{x})_i$. In total, the output covariability
224 \mathbf{C} will simply inherit the rank of the input covariability \mathbf{X} . Attentional modulation affects c_i through
225 \mathbf{K} and \mathbf{L} and we easily satisfy condition **C1** with $g_i = c_i^A / c_i^U$.

226 What remains is to find constraints on \mathbf{J} and the attentional modulation of \mathbf{L} that satisfy condition
227 **C2**. Let us consider the case where $c_i^U, c_i^A > 0$ so that condition **C2** is satisfied when $c_i^A - c_i^U < 0$.
228 For the sake of mathematical simplicity let us separate the population into qN excitatory neurons and
229 $(1 - q)N$ inhibitory neurons ($0 < q < 1$). Let all excitatory (inhibitory) neurons project with synaptic

strength J_E ($-J_I$), have gain L_E (L_I), and receive the external inputs of strength x_E (x_I). Finally, let the probability for all connections be p , and consider only weak connections ($J \propto 1/N$ and N large) so that we can ignore the influence of polysynaptic paths in the network (Pernice et al., 2011; Trousdale et al., 2012). Then the attentional modulation of an excitatory neuron decomposes into:

$$c_E^A - c_E^U = \underbrace{\left(L_E^A - L_E^U \right) x_E}_{\text{direct external input}} + \underbrace{\left(L_E^A - L_E^U \right) q p N J_E x_E}_{\text{external input filtered through the excitatory population}} - \underbrace{\left(L_I^A - L_I^U \right) (1-q) p N J_I x_I}_{\text{external input filtered through the inhibitory population}}. \quad (5)$$

The first term is the direct transfer of the external fluctuations, and the second and third terms are indirect transfer of external fluctuations via the excitatory and inhibitory populations, respectively. Recall that $L^A - L^U > 0$, meaning that for $c_E^A - c_E^U < 0$ to be satisfied we require the third term to outweigh the combination of the first and second terms. In other words, the inhibitory population must experience a sizable attentional modulation. A similar cancelation of correlations by recurrent inhibition has been recently studied in a variety of cortical models (Renart et al., 2010; Tetzlaff et al., 2012; Ly et al., 2012; Doiron et al., 2016; Rosenbaum et al., 2017).

In the above we considered weak synaptic connections where $J_{ij} \sim 1/N$. Rather, if we scale $J_{ij} \sim 1/\sqrt{N}$, as would be the case for classical balanced networks (van Vreeswijk and Sompolinsky, 1998), then for very large N the solution no longer depends upon the gain L . Finite N or the inclusion of synaptic nonlinearities through short term plasticity (Mongillo et al., 2012) may be necessary to satisfy condition **C2** with large synapses. Furthermore, the large synaptic weights associated with $J_{ij} \sim 1/\sqrt{N}$ do not allow us to neglect polysynaptic paths, as is needed for Eq. (5). Extending our analysis to networks with balanced scaling will be the focus of future work.

In summary our analysis has identified two circuit features that allow recurrent networks to capture conditions **C1** and **C2** for attentional modulation. First, the network must be subject to a global source of external fluctuations that dominates network covariability (**C1**). Second, the network must have recurrent inhibitory connections that are subject to a large attentional modulation (**C2**).

252 Mean field model of attention

We next apply the intuition gained in the preceding section to propose a cortical model that captures key neural correlates of attentional modulation. We model V4 as a recurrently coupled network of excitatory and inhibitory leaky integrate-and-fire model neurons (Tetzlaff et al., 2012; Ledoux and Brunel, 2011; Trousdale et al., 2012; Doiron et al., 2004) (Fig. 4a). In addition to recurrent synaptic inputs, each neuron receives private and global sources of external fluctuating input (Fig. 4b). The global noise is an attention-independent source of input correlation that the network filters and transforms into network-wide output spiking correlations (Fig. 4c).

While the linear response theory introduced in Eq. (3) is well suited to study large networks of integrate-and-fire neurons driven by weakly correlated inputs (Tetzlaff et al., 2012; Ledoux and Brunel, 2011; Trousdale et al., 2012; Doiron et al., 2004), the analysis offers little analytic insight. Instead, we consider the instantaneous activity across population α : $r_\alpha(t) = \frac{1}{N_\alpha} \sum_i y_{i\alpha}(t)$, where $y_{i\alpha}(t)$ is the spike train from neuron i of population α and N_α is the population size ($\alpha = E$ or I). This approach reduces the model to just the two dynamic variables, the excitatory population rate $r_E(t)$ and the inhibitory population rate $r_I(t)$ ($r_E(t)$ is shown in Fig. 4d). Despite this severe reduction the model retains the key ingredients for attentional modulation identified in the previous section – recurrent excitation and inhibition combined with a source of global fluctuations.

We take the population sizes to be large and consider a phenomenological dynamic mean field (Tetzlaff et al., 2012; Ledoux and Brunel, 2011) of the cortical network (see Methods: Mean field model):

$$\begin{aligned} \tau_E \frac{dr_E}{dt} &= -r_E + f_E (\mu_E + J_{EE} r_E - J_{EI} r_I + \sigma_E \xi(t)), \\ \tau_I \frac{dr_I}{dt} &= -r_I + f_I (\mu_I + J_{IE} r_E - J_{II} r_I + \sigma_I \xi(t)). \end{aligned} \quad (6)$$

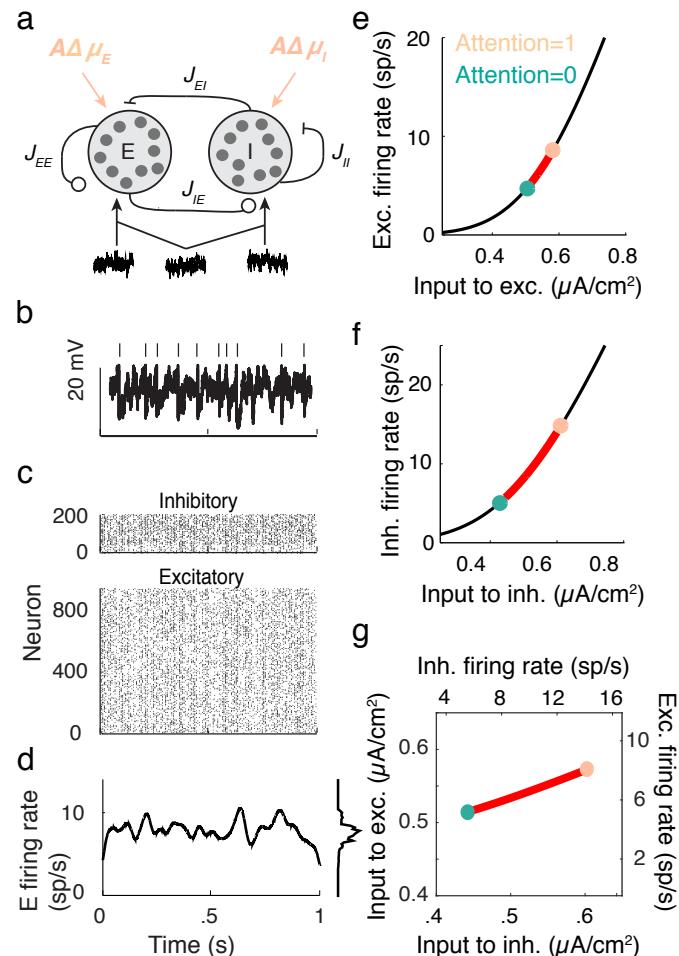


Figure 4. Excitatory-inhibitory network model. **a**, Recurrent excitatory-inhibitory network subject to private and shared fluctuations as well as top-down attentional modulation. **b**, Example voltage trace from a LIF model neuron in the network. Top tick marks denote spike times. **c**, Spike time raster plot of the spiking activity from the model network. **d**, Population-averaged firing rate $r_E(t)$ of the excitatory population. Left: frequency distribution of population-averaged firing rate. **e**, Transfer function f_E between the effective input and the firing rate for a model excitatory neuron. The red segment represents the attentional shift in effective input and hence firing rate. **f**, Same as e, but for the inhibitory population. **g**, Attention as a path through (\bar{r}_E, \bar{r}_I) space, and equivalently through $(I_E^{\text{eff}}, I_I^{\text{eff}})$ space.

269 The function f_α is the input-output transfer of population α , taken to be the mean firing rate for a
 270 fixed input (Figs. 4e for the E population and 4f for the I population). The parameter $J_{\alpha\beta}$ is the
 271 coupling strength from population β to population α . Finally, μ_α and σ_α are the respective strengths
 272 of the mean input and the global fluctuation $\xi(t)$ to population α (throughout $\xi(t)$ has a zero mean).
 273 To simplify our exposition we take symmetric coupling $J_{EE} = J_{IE} \equiv J_E$ and $J_{EI} = J_{II} \equiv J_I$ and
 274 symmetric timescales $\tau_E = \tau_I (= 1)$. We set the recurrent coupling so that the model has a stationary
 275 mean firing rate (\bar{r}_E, \bar{r}_I) , about which $\xi(t)$ induces fluctuations in $r_E(t)$ and $r_I(t)$.

276 Attention is modeled as a top-down influence on the static input: $\mu_\alpha = \mu_{\alpha B} + A\Delta\mu_\alpha$. Here $\mu_{\alpha B}$ is
 277 a background input, the parameter A models attention with $A = 0$ denoting the unattended state
 278 and $A = 1$ the fully attended state, and $\Delta\mu_\alpha > 0$ is the increase in μ_α due to attention. We note
 279 that the choice of representing the unattended state by $A = 0$ and the attended state by $A = 1$ is
 280 only due to convenience, and is not meant to make any statement about particular bounds on these
 281 states. In this model attention simply increases the excitability of all of the neurons in the network
 282 (Fig. 4a). This modulation is consistent with the rank one structure of attentional modulation in
 283 the data (Fig. 2), since μ_α is a single neuron property. The attention-induced increase in (μ_E, μ_I)
 284 causes an increase in the mean firing rates (\bar{r}_E, \bar{r}_I) (red paths in Figs. 4e,f), consistent with recordings
 285 from putative excitatory (*McAdams and Maunsell, 2000; Reynolds et al., 1999*) and inhibitory
 286 neurons (*Mitchell et al., 2007*) in visual area V4. Since f_α is a simple rising function then there is
 287 a unique mapping of an attentional path in (μ_E, μ_I) space to a path in (\bar{r}_E, \bar{r}_I) space (Fig. 4g).

288 In total, our population model has the core features required to satisfy Conditions **C1** and **C2**
 289 of the previous section. We next use our mean field model to investigate how attentional paths in
 290 (\bar{r}_E, \bar{r}_I) space affect population spiking variability.

291 **Attention modulates population variability**

292 The global input $\xi(t)$ causes fluctuations about the network stationary state: $r_\alpha(t) = \bar{r}_\alpha + \delta r_\alpha(t)$. The
 293 fluctuations $\delta r_\alpha(t)$ are directly related to coordinated spiking activity in population α . In particular,
 294 in the limit of large N_α we have that $V_E \equiv \text{Var}(r_E) \propto \langle \text{Cov}(y_i, y_j) \rangle$, where the expectation is over
 295 (i, j) pairs in the spiking network. Thus, in our mean field network we require attentional modulation
 296 to decrease population variance V_E .

297 For sufficiently small σ_α the fluctuations $\delta r_E(t)$ and $\delta r_I(t)$ obey linearized mean field equations
 298 (see Methods: Mean field model, Eq. (17)). The linear system is readily analyzed and we obtain the
 299 population variance V_E computed over long time windows (see Methods: Computing V_E):

$$V_E = \left[\frac{L_E(J_I L_I(\sigma_E - \sigma_I) + \sigma_E)}{1 + J_I L_I - J_E L_E} \right]^2. \quad (7)$$

300 Here $L_\alpha \equiv f'_\alpha$ is the response gain of neurons in population α . Equation (7) shows that V_E depends
 301 directly on L_α , and we recall that L_α changes with attention (the slope of f_α in Fig. 4e,f). Thus,
 302 while the derivation of V_E requires linear fluctuations about a steady state, attentional modulation
 303 samples the nonlinearity in the transfer f_α by changing the state about which we linearize. Any
 304 attention-mediated change in V_E is not obvious since both $L_I^A > L_I^U$ and $L_E^A > L_E^U$, meaning that
 305 both the numerator and denominator in Eq. (7) will change with attention.

306 We explore V_E by sweeping over (\bar{r}_E, \bar{r}_I) space (Fig. 5a). When the network has high \bar{r}_E and low
 307 \bar{r}_I then V_E is large, while V_E is low for the opposite case of high \bar{r}_I and low \bar{r}_E . Along our attention
 308 path r_E increases while V_E decreases (Fig. 5b), satisfying our requirements for attentional modulation.
 309 The attention path that we highlight is just one potential path that reduces population variability,
 310 however all paths which reduce V_E share a large attention-mediated recruitment of inhibition. If we
 311 start with the unattended state (turquoise dot in Fig. 5c) we can label all $(\Delta\mu_E > 0, \Delta\mu_I > 0)$ points
 312 that have a smaller population variance than the unattended point (light green region in Fig. 5c).
 313 These modulations all share that $\Delta\mu_I > \Delta\mu_E$ (Fig. 5c, green region is below the $\Delta\mu_E = \Delta\mu_I$ line).
 314 While the absolute comparison between $\Delta\mu_E$ and $\Delta\mu_I$ may depend on model parameters, a robust
 315 necessary feature of top-down attentional modulation is that it must significantly recruit the inhibitory
 316 population. This observation is a major circuit prediction of our model.

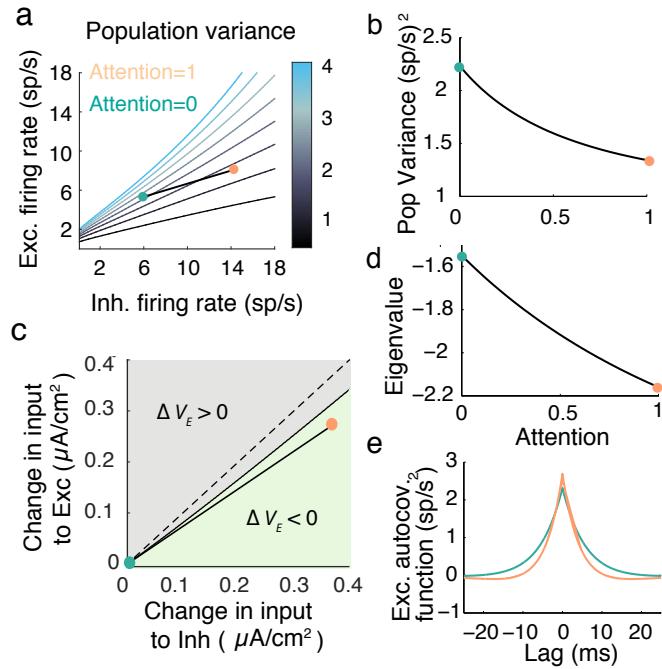


Figure 5. Mean field model shows an attention mediated decrease in population variance. **a**, An attentional path in excitatory-inhibitory firing rate space for which the population variance decreases. Colored contours define iso-lines of population variance in increments of 10 (sp/s^2). The attentional path links the unattended state ($A = 0$; turquoise point) to the attended state ($A = 1$, orange point). **b**, Variance values as a function of the attentional path defined in **a**. **c**, The modulation from an unattended state (origin) to an attended state over the input space ($\Delta\mu_E, \Delta\mu_I$). Solid black line marks where V_E remains unchanged, and the green region where $\Delta V_E = \text{Var}^A(r_E) - \text{Var}^U(r_E)$ is less than zero. **d**, The eigenvalue (λ) along the attentional path. With increased attention it becomes more negative, indicating that the state (\bar{r}_E, \bar{r}_I) is more stable. **e**, Autocovariance function of the excitatory population rate $r_E(t)$ in the attended and unattended state (computed using Eq. (19)).

317 An intuitive way to understand inhibition's role in the decrease in population variance is through
 318 the stability analysis of the mean field equations. The eigenvalues of the linearized system are
 319 $\lambda_1 = -1 - J_I L_I + J_E L_E < 0$ and $\lambda_2 = -1$ (see Methods: Mean field model, Eq. (18)). Note that
 320 the denominator of the population variance (Eq. 7) equals the square of the eigenvalue product
 321 $\lambda_1 \lambda_2 = 1 + J_I L_I - J_E L_E$. The stability of the network activity is determined by λ_1 ; the more negative
 322 λ_1 , the more stable the point (\bar{r}_E, \bar{r}_I) , and the better the network dampens the perturbations about
 323 the point due to input fluctuations $\xi(t)$. The decrease of λ_1 along the example attention path is
 324 clear (Fig. 5d), and overcomes the increase in the numerator of V_E due to increases in L_E and L_I .
 325 The enhanced damping is why V_E decreases, explicitly seen in the steeper decline of the excitatory
 326 population autocovariance function in the attended compared to the unattended state (Fig. 5e).

327 This enhanced stability due to recurrent inhibition is a reflection of inhibition canceling population
 328 variability provided by external fluctuations and recurrent excitation (Renart *et al.*, 2010; Tetzlaff
 329 *et al.*, 2012; Ozeki *et al.*, 2009). Indeed, taking the coupling J to be weak allows the expansion
 330 $(1 + J_I L_I - J_E L_E)^{-2} \approx 1 + 2J_E L_E - 2J_I L_I$ in Eq. (7), so that the attention mediated increase in
 331 L_I reduces population variance through cancellation, as in Eq. (5). However, this expansion is not
 332 formally required to compute the eigenvalues λ_1 and λ_2 , and these measure the stability of the firing
 333 rate dynamics. We mention the expansion only to compare to the original motivation for inhibition.

334 The expression for V_E given above (Eq. 7) assumes a symmetry in the network coupling, namely
 335 that $J_{EE} = J_{IE} \equiv J_E$ and $J_{EI} = J_{II} \equiv J_I$. This allowed V_E to be compactly written, facilitating
 336 the analysis of how attention affects both the numerator and denominator of Eq. (7). However, the

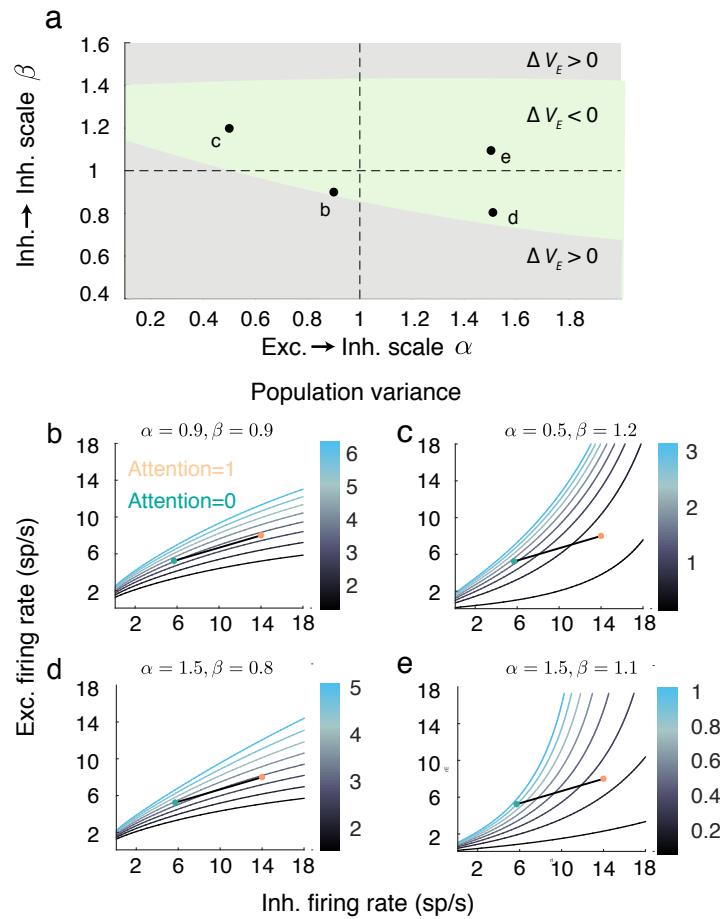


Figure 6. The attention mediated reduction in population variance is robust to changes in strength of recurrent connectivity. **a**, Sweep over $\alpha = J_{EE}/J_{IE}$ and $\beta = J_{EI}/J_{II}$ space (with J_{EE} and J_{EI} fixed) labeling the region where $\Delta V_E = V_E^U - V_E^A$ is positive (grey) and negative (green). **b-e**, Attentional path in excitatory-inhibitory firing rate space. The colored contours are as in Fig. 5a. All calculations are done using Eqs. (18)-(20).

linearization of the mean field equations and the subsequent analysis of population variability do not require this assumption (see Methods: Mean field model Eqs. (18)-(20)). To explore the robustness of our main result we let $J_{IE} = \alpha J_E$ and $J_{II} = \beta J_I$, thereby breaking the coupling symmetry for $\alpha, \beta \neq 1$. The reduction in V_E with attention is robust over a large region of (α, β) (Fig. 6a, green region). Focusing on selected (α, β) pairings within the region where V_E decreases shows that the attentional path identified for the network with coupling symmetry produces qualitatively similar behavior in the more general network (compare Fig. 5c to Fig. 6b-e). In total, the inhibitory mechanism for attention mediated reduction in population variability is robust to changes in the recurrent coupling with the network.

While the reduced mean field equations are straightforward to analyze, a similar attenuation of pairwise covariance $\text{Cov}(y_i, y_j)$ along the same attentional path occurs in the LIF model network (Appendix: Spiking network). Using linear response analysis for the spiking network we can relate the effect of inhibition to previous work in spiking networks (Renart et al., 2010; Tetzlaff et al., 2012; Ly et al., 2012; Doiron et al., 2016). In particular, the attention-mediated decrease of $\text{Cov}(y_i, y_j)$ occurs for a wide range of timescale, ranging as low as 20ms. However, for short timescales that match the higher gamma frequency range (approximately 60-70 Hz) this attentional modulation increases $\text{Cov}(y_i, y_j)$ (Appendix Fig. 6). This finding is consistent with reports of attention-mediated increases

354 of neuronal synchrony on gamma frequency timescales (*Fries et al., 2001; Buia and Tiesinga,*
355 *2008*), particularly when inhibitory circuits are engaged (*Kim et al., 2016*).

356 **Attention can simultaneously increase stimulus gain and decrease noise covariance**
357 An important neural correlate of attention is enhanced stimulus response gain (*McAdams and*
358 *Maunsell, 2000*). The previous section outlines how the recruitment of recurrent inhibitory feedback
359 by attention reduces response variability. However, inhibitory feedback is also a common gain control
360 mechanism, and increased inhibition reduces response gain through the same mechanism that dampens
361 population variability (*Sutherland et al., 2009*). Thus it is possible that the decorrelating effect of
362 attention in our model may also reduce stimulus response gain as well, which would make the model
363 inconsistent with experimental data.

364 To insert a bottom-up stimulus s in our model we let the attention-independent background input
365 have a stimulus term: $\mu_{\alpha B} = k_{\alpha} s + \hat{\mu}_{\alpha B}$. Here k_{α} is the feedforward stimulus gain to population
366 α and $\hat{\mu}_{\alpha B}$ is the background input that is both attention and stimulus independent. Our model
367 captures a bulk firing rate r_E rather than a population model with distributed tuning. Because of this
368 the stimulus s should either be conceived as the contrast of an input, or the population conceived as a
369 collection of identically-tuned neurons (i.e a single cortical column).

370 Straightforward analysis shows that the stimulus response gain of the excitatory population can be
371 written as (Methods: Computing stimulus response gain):

$$G_E \equiv \frac{d\bar{r}_E}{ds} = \frac{k_E \sqrt{V_E}}{\sigma_E} + \frac{J_I L_E L_I}{1 + J_I L_I - J_E L_E} (k_E - k_I). \quad (8)$$

372 If $k_E = k_I$ then $G_E \propto \sqrt{V_E}$, and thus any attentional modulation that reduces population variability
373 will necessarily reduce population stimulus sensitivity. However, for $k_E > k_I$ the second term in Eq.
374 (8) can counteract this effect and decouple stimulus sensitivity and variability modulations.

375 Consider the example attentional path (Fig. 4g) with the extreme choice of $k_E = 1$ and $k_I = 0$. In
376 this case attention causes an increase in G_E (Fig. 7a,b), while simultaneously causing a decrease in
377 V_E (Fig. 5a,b). This is a robust effect, as seen by the region in (\bar{r}_E, \bar{r}_I) space for which the change in
378 V_E from the unattended state is negative, and the change in G_E is positive (green region, Fig. 7c).
379 Further, for fixed k_I the proportion of the gray rectangle that the green region occupies increases with
380 $k_E > k_I$ (Fig. 7d). Thus, the decoupling of attentional effects on population variability and stimulus
381 sensitivity is robust to both attentional path ($\Delta\mu_E, \Delta\mu_I$) and feedforward gain (k_E, k_I) choices. The
382 condition that $k_E > k_I$ implies that feedforward stimuli must directly target excitatory neurons to
383 a larger degree than inhibitory neurons (or at least the inhibitory neurons subject to attentional
384 modulation). This gives us a complementary prediction to the one from the previous section: while
385 top-down attention favors inhibitory neurons, the bottom-up stimulus favors excitatory neurons.

386 In total, our model of attentional modulation in recurrently coupled excitatory and inhibitory
387 cortical networks subject to global fluctuations satisfies three main neural correlates of attention: (1)
388 increase in excitatory firing rates and in (2) stimulus-response gain, with a (3) decrease in pairwise
389 excitatory neuron co-variability.

390 Impact of attentional modulation on neural coding

391 Attention serves to enhance cognitive performance, especially on discrimination tasks that are difficult
392 (*Moore and Zirnsak, 2017*). Thus, it is expected that the attention-mediated reduction in popu-
393 lation variability and increase in stimulus response gain subserve an enhanced stimulus estimation
394 (*Cohen and Maunsell, 2009; Ruff and Cohen, 2014*). In this section we investigate how the
395 attentional modulation outlined in the previous sections affects stimulus coding by the population.

396 As mentioned above our simplified mean field model (Eq. 6) considers only a bulk response, where
397 any individual neuron tuning is lost. As such a proper analysis of population coding is not possible.
398 Nonetheless, our model has two basic features often associated with enhanced coding, decreased
399 population variability (Fig. 5) and increased stimulus-response gain (Fig. 7).

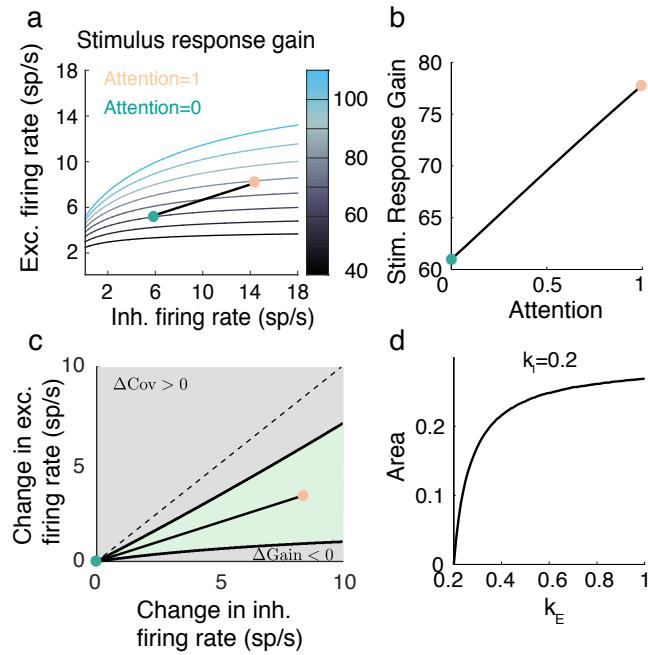


Figure 7. Attention model can capture increase in stimulus response gain G_E despite decrease in population variance V_E . **a**, Attentional path through (\bar{r}_E, \bar{r}_I) space shows an increase in stimulus response gain. The shown path is the same path as in Figure 5. **b**, Values of G_E along the path depicted in a). **c**, The green region in (\bar{r}_E, \bar{r}_I) space denotes where $\Delta V_E = \text{Var}^A(r_E) - \text{Var}^U(r_E) < 0$ and $\Delta G_E = G_E^A - G_E^U > 0$. Black lines are iso-lines of covariance and gain, along which those quantities do not change. **d**, Percent area of the green region in c) out of a constant rectangle, as the feedforward stimulus gain k_E increases, with $k_I = 0.2$ held constant.

Fisher information (Averbeck *et al.*, 2006; Beck *et al.*, 2011) gives a lower bound on the variance of a stimulus estimate constructed from noisy population responses, and is an often used metric for population coding. The linear Fisher information (Beck *et al.*, 2011) FI_{EI} computed from our two-dimensional recurrent network is:

$$\text{FI}_{EI} = \begin{bmatrix} G_E & G_I \\ C_{EI} & V_I \end{bmatrix} \begin{bmatrix} V_E & C_{EI} \\ C_{EI} & V_I \end{bmatrix}^{-1} \begin{bmatrix} G_E \\ G_I \end{bmatrix} = \text{constant} \quad (9)$$

400 Here $V_\alpha = \text{Var}(r_\alpha)$, $G_\alpha = d\bar{r}_\alpha/ds$, and $C_{EI} = \text{Cov}(r_E, r_I)$. The important result is that FI_{EI} is
401 invariant with attention, meaning that attention does not increase the network's capacity to estimate
402 the stimulus s .

403 While the proof of Eq. (9) is straightforward and applies to our recurrent excitatory-inhibitory
404 population (see Methods: Fisher information), the invariance of the total information F_{EI} with attention
405 is most easily understood by analogy with an uncoupled, one-dimensional excitatory population (Fig.
406 8a). Without coupling, the input to the population is simply $k_E s + \sigma_E \xi(t)$, which is then passed
407 through the firing rate nonlinearity f_E . In this case the gain is $G_E = k_E L_E$, and assuming a linear
408 transfer the population variance is $V_E = \sigma_E^2 L_E^2$. In total the linear Fisher information from the
409 uncoupled population is then:

$$\text{FI}_E^{\text{uc}} = \frac{G_E^2}{V_E} = \frac{(k_E L_E)^2}{\sigma_E^2 L_E^2} = \frac{k_E^2}{\sigma_E^2}. \quad (10)$$

410 The proportion L_E^2 by which attention increases the squared gain (Fig. 8a, top) is exactly matched by
411 the attention related increase in population variance (Fig. 8a, bottom), resulting in cancellation of
412 any attention-dependent terms in FI_{EI} .

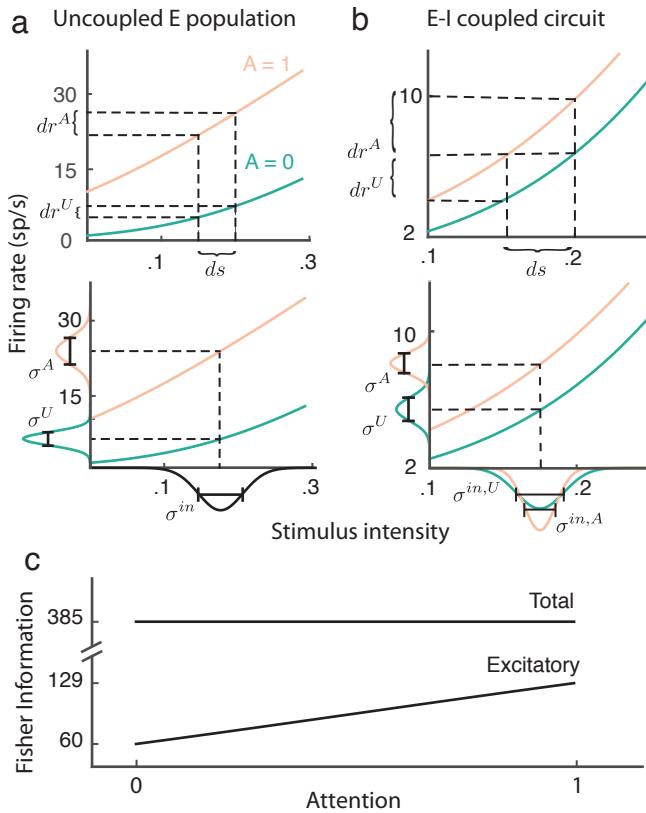


Figure 8. Attention improves stimulus estimation by the excitatory population embedded within excitatory (E)-inhibitory (I) network. **a**, Top: For a uncoupled excitatory population, the stimulus response gain G_E increases with attention. Turquoise: unattended state; orange: attended state. Bottom: Population variance V_E increases with attention. Stimulus-response curves same as above. Input variance is computed from all input to a population, including external noise and recurrent coupling. The Fisher information for the uncoupled E population is constant with attention because the squared gain G_E^2 and variance V_E increase proportionally **b**, Same as a) but for the E population within the $E - I$ network. Top: G_E increases with attention. Bottom: V_E decreases with attention, because the net input variance of the E population decreases with attention. **c**, Total Fisher information for coupled E-I populations is constant with attention. By contrast, the Fisher information of the excitatory component FI_E increases with attention.

413 The majority of projection neurons in the neocortex are excitatory, so we now consider the stimulus
 414 estimation from a readout of only the excitatory population. Combining our previous results we obtain:

$$415 \quad FI_E = \frac{G_E^2}{V_E} = \frac{(J_I L_I (k_E - k_I) + k_E)^2}{\sigma_E^2 - J_I^2 L_I^2 (\sigma_E \sigma_I - \sigma_E^2 - \sigma_I^2) - 2 J_I L_I \sigma_E (\sigma_I - \sigma_E)}. \quad (11)$$

416 Restricting the readout to be from only the excitatory population drastically reduces the total
 417 information (compare FI_{EI} to FI_E in Fig. 8c). As with the uncoupled population the response gain
 418 G_E of the excitatory neurons in the coupled population increases with attention (Fig. 8b, top). Yet
 419 unlike the uncoupled population the net input variability to the E population is reduced by attention
 420 through a cancelation of the external variability $\xi(t)$ via inhibition (Fig. 8b, bottom). These two
 421 components combine so that despite $FI_E < FI_{EI}$, we have that FI_E does increase with attention (Fig.
 422 8c). In sum, even though the total stimulus information in the network does not change with attention,
 423 the amount of information extractable from the excitatory population increases, which could lead to
 424 improved downstream stimulus estimation in the attended state.

425 Discussion

426 Using population recordings from visual area V4 we identified rank one structure in the mapping of
427 population spike count covariability between unattended and attended states. We used this finding to
428 motivate an excitatory-inhibitory cortical circuit model that captures both the attention-mediated
429 increases in the firing rate and stimulus response gain, as well as decreases in noise correlations. Our
430 model accomplishes this with only an attention dependent shift in the overall excitability of the cortical
431 population, in contrast to a scheme where distinct biophysical mechanisms would be responsible for
432 respective firing rate and noise correlations modulations. The model makes two key predictions about
433 how stimulus and modulatory inputs are distributed over the excitatory-inhibitory cortical circuit.
434 First, top-down attentional signals must affect inhibitory neurons more than excitatory neurons to
435 allow a better damping of global fluctuations in the attended state. Second, bottom-up stimulus
436 information must be biased towards excitatory cells to permit higher gain in the attended state. In
437 total, the increased response gain and decreased correlations enhance the flow of information when the
438 readout is confined to the excitatory population.

439 Candidate physiological mechanisms for attentional modulation

440 Our model does not consider a specific type of inhibitory neuron, and rather models a generic recurrent
441 excitatory-inhibitory circuit. However, inhibitory circuits in cortex are complex, with at least three
442 distinct interneuron types being prominent in many areas: parvalbumin- (PV), somatostatin- (SOM),
443 and vasointestinal peptide-expressing (VIP) interneurons (*Rudy et al., 2011; Pfeffer et al., 2013;*
444 *Kepecs and Fishell, 2014*). In mouse visual cortex, both SOM and PV cells form recurrent circuits
445 with pyramidal cells, with PV cells having stronger inhibitory projections to pyramidal cells than those
446 of SOM cells (*Pfeffer et al., 2013*). Furthermore, PV and SOM neurons directly inhibit one another,
447 with the SOM to PV connection being stronger than the PV to SOM connection (*Pfeffer et al., 2013*,
448 *2013*). Finally, VIP cells project strongly to SOM cells (*Pfeffer et al., 2013*) and are activated from
449 inputs outside of the circuit (*Lee et al., 2013; Fu et al., 2014*), making them an attractive target for
450 modulation. Recent studies in visual, auditory, and somatosensory cortical circuits show that VIP cell
451 activation provides an active disinhibition of pyramidal cells via a suppression of SOM cells (*Kepecs*
452 *and Fishell, 2014*). Basal forebrain (BF) stimulation modulates both muscarinic and nicotinic ACh
453 receptors (mAChRs and nAChRs respectively) in a fashion that mimics attentional modulation (*Alitto*
454 *and Dan, 2012*). In particular, the recruitment of VIP cell activity *in vivo* through BF stimulation
455 is strongly dependent on both the muscarinic and nicotinic cholinergic pathways (*Alitto and Dan, 2012;*
456 *Kuchibhotla et al., 2017; Fu et al., 2014*), and it has thus been hypothesized VIP cells
457 activation could be an important component of attentional modulation (*Alitto and Dan, 2012;*
458 *Poorthuis et al., 2014*).

459 If we consider the inhibitory population in our model to be PV interneurons then the recruitment
460 of VIP cell activity via top-down cholinergic pathways is consistent with our attentional model in
461 two ways. First, activation of the VIP → SOM → pyramidal cell pathway provides a disinhibition
462 to pyramidal cells, modeled simply as an overall depolarization to pyramidal cells in the attended
463 state (Fig. 4). Second, the activation of the VIP → SOM → PV cell pathway disinhibits PV cells,
464 and the strong SOM → PV projection would suggest that the disinhibition is sizable as required by
465 our model (Fig. 5c). Finally, a recent study in mouse medial prefrontal cortex reports that identified
466 PV interneurons show an attention related increase in activity, and that optogenetic silencing of PV
467 neurons impairs attentional processing (*Kim et al., 2016*).

468 However, our logic is perhaps overly simplistic and neglects the direct modulation of SOM cells
469 via muscarinic and nicotinic cholinergic pathways (*Alitto and Dan, 2012; Kuchibhotla et al.,*
470 *2017*) that could compromise the disinhibitory pathways. Further, there is evidence of a direct ACh
471 modulation of PV cells (*Disney et al., 2014*) as opposed to through a disinhibitory pathway. Finally,
472 there may be important differences across both species (mouse vs. primate) and visual area (V1 vs.
473 V4) that fundamentally change the pyramidal, PV, SOM, and VIP circuit that is understood from
474 mouse V1 (*Pfeffer et al., 2013*). Future studies in the inhibitory to excitatory circuitry of primate

475 visual cortex, and its attentional modulation via neuromodulation, are required to navigate these
476 issues.

477 Finally, the simultaneous increase in response gain and decrease in noise correlations with attention
478 requires excitatory neurons to be more sensitive to bottom-up visual stimulus than inhibitory neurons
479 ($k_E > k_I$, Fig. 7). In mouse visual cortex, GABAergic interneurons show overall less stimulus selectivity
480 than pyramidal neurons (Sohya *et al.*, 2007), however this involves both direct feedforward and
481 recurrent contributions to stimulus tuning. While our model simplified the feedforward stimulus
482 gain k_E and k_I to be constant with attention, it is known that attention also modulates feedforward
483 gain through presynaptic nACh receptors (Disney *et al.*, 2007). Notably, nAChRs are found at
484 thalamocortical synapses onto layer 4 excitatory cells and not onto inhibitory neurons, suggesting that
485 k_E would increase with attention while k_I would not. Thus, k_E should also increase with attention
486 while k_I should not, further supporting that $k_E > k_I$.

487 Modeling global network fluctuations and their modulation

488 Our model considered the source of global fluctuations as external to the network. This choice was
489 due in part to difficulties in producing global, long timescale fluctuations through strictly internal
490 coupling (Renart *et al.*, 2010; Rosenbaum *et al.*, 2017). Our model assumed that the intensity of
491 these external input fluctuation were independent of attention. Rather, attention shifted the operating
492 point of the network such that the transfer of input variability to population-wide output activity was
493 attenuated in the attended state.

494 Recent analysis of population recordings show that generative models of spike trains that consider
495 gain fluctuations in conjunction with standard spike emission variability capture much of the variability
496 of cortical dynamics (Rabinowitz *et al.*, 2015; Lin *et al.*, 2015). Further, these gain fluctuations
497 are well approximated by a one-dimensional, global stochastic process affecting all neurons in the
498 population (Ecker *et al.*, 2014; Rabinowitz *et al.*, 2015; Lin *et al.*, 2015; Ecker *et al.*, 2016;
499 Engel *et al.*, 2016; Whiteway and Butts, 2017). When these techniques are applied to population
500 recordings subject to attentional modulation, the global gain fluctuations are considerably reduced in
501 the attended state (Rabinowitz *et al.*, 2015; Ecker *et al.*, 2016). Our assumption that external
502 input fluctuations to our network are attention-invariant is consistent with this statistical analysis since
503 it is necessarily constructed from only output activity. Nevertheless, another potential model is that
504 the reduction in population variability is simply inherited from an attention-mediated suppression of
505 the global input fluctuations. Unfortunately, it is difficult to distinguish between these two mechanisms
506 when restricted to only output spiking activity.

507 However, a model where output variability reductions are simply inherited from external inputs
508 suffers from two criticisms. First, it begs the question: what is the mechanism behind the shift in input
509 variability? Second, our model requires only an increase in the external depolarization to excitatory
510 and inhibitory populations to account for all attentional correlates. An inheritance model would
511 necessarily decouple the attentional mechanisms behind increases in network firing rate (still requiring
512 a depolarization) and the decrease in global input variability. Thus, our model offers a parsimonious
513 and biologically motivated explanation of these neural correlates of attention. Further work dissecting
514 the various external and internal sources of variability to cortical networks, and their attentional
515 modulation, is needed to properly validate or refute these different models.

516 Attentional modulation of neural coding through inhibition

517 Our network model assumed attention-invariant external fluctuations and weak recurrent inputs,
518 permitting a linear analysis of network activity. As a consequence the linear information transfer by
519 the entire population was attention-invariant (Fig. 8), because attention modulated the network's
520 transfer of signal and noise equivalently. However, this invariance was only apparent if the decoder
521 had access to both the excitatory and inhibitory populations. However, most of the neurons in cortex
522 that project between areas are excitatory. When the decoder was restricted to only the activity of the
523 excitatory population then our analysis uncovered two main results. First, the excitatory population

524 carried less information than the combined excitatory-inhibitory activity, suggesting an inherently
 525 suboptimal coding scheme used by the cortex. Second, the attention-mediated modulation of the
 526 inhibitory neurons increased the information carried by the excitatory population. This agrees with the
 527 wealth of studies that show that attention improves behavioral performance on stimulus discrimination
 528 tasks.

529 Determining the impact of population-wide spiking variability on neural coding is complicated (*Aver-*
 530 *beck et al., 2006; Kohn et al., 2016*). A recent theoretical study has shown that noise correlations
 531 that limit stimulus information must be parallel to the direction in which population activity encodes
 532 the stimulus (*Moreno-Bote et al., 2014*). The fluctuations in our network satisfy this criteria,
 533 albeit trivially since all neurons share the same stimulus input. Indeed, in our network the external
 534 inputs appear to the network as $s + x(t)$, meaning that fluctuations from the noise source $x(t)$ are
 535 indistinguishable from fluctuations in the stimulus s . This is an oversimplified view and assumes that
 536 the decoder treats the neurons as indistinguishable from one another, at odds with classic work in
 537 population coding (*Pouget et al., 2000*). Extending our network to include distributed tuning and
 538 feature-based recurrent connectivity is a natural next step (*Ben-Yishai et al., 1995; Rubin et al.,*
 539 *2015*). To do this the spatial scales of feedforward tuning, recurrent projections, external fluctuations,
 540 as well as attention modulation must all be specified. It is not clear how noise correlations will depend
 541 on these choices yet work in spatially distributed balanced networks shows that solutions can be
 542 complex (*Rosenbaum et al., 2017*).

543 The role of inhibition in shaping cortical function is a longstanding topic of study (*Isaacson*
 544 *and Scanziani, 2011*), including recent work showing inhibition can actively decorrelate cortical
 545 responses (*Renart et al., 2010; Tetzlaff et al., 2012; Ly et al., 2012*). Our work gives a concrete
 546 example of how this decorrelation can be gated and used to control the flow of information. A natural
 547 extension of this work is to consider neural populations with distributed tuning, so that each cell has a
 548 distinct preferred stimulus. Importantly, in tasks that probe a distributed population attention again
 549 decreases noise correlations between neurons with similar stimulus preference, yet it *increases* noise
 550 correlations between cells with dissimilar stimulus preference (*Ruff and Cohen, 2014*). The circuit
 551 mechanisms underlying this neural correlate of attention are unclear. However, there is ample work in
 552 understanding how recurrent inhibition shapes cortical activity in distributed populations (*Isaacson*
 553 *and Scanziani, 2011*), including in models of attentional circuits (*Aridid et al., 2007; Buia and*
 554 *Tiesinga, 2008*). Adapting our model to include distributed tuning is an important next step and
 555 will be a better framework to discuss the coding consequences of the attentional modulation circuits
 556 proposed in our study.

557 Methods and Materials

558 Data preparation

559 Data was collected by from two rhesus monkeys with microelectrode arrays implanted bilaterally in V4
 560 as they performed an orientation-change detection task (Figure 1a) (*Cohen and Maunsell, 2009*).
 561 All animal procedures were in accordance with the Institutional Animal Care and Use Committee of
 562 Harvard Medical School. Two oriented Gabor stimuli flashed on and off several times, until one of
 563 them changed orientation. The task of the monkey was to then saccade to the stimulus that changed.
 564 Each recording session consisted of at least four blocks of trials in which the monkey's attention was
 565 cued to the left or right. We excluded from the analysis instruction trials which occurred at the start of
 566 each block to cue the monkey to one side to attend to, catch trials in which the monkey was rewarded
 567 just for fixating, and trials in which the monkey did not perform the task correctly. Moreover, the
 568 first and last stimulus presentations in each trial were not analyzed, to prevent transients due to
 569 stimulus appearance or change from affecting the results. The total number of trials included in the
 570 analysis from all the recording sessions was 42,496. Each trial consisted of between 3 and 12 stimulus
 571 presentations, of which all but the first and last were analyzed.

572 Recordings from the left and right hemispheres of each monkey were analyzed separately because the
 573 activities of the neurons in opposite hemispheres had near-zero correlations (*Cohen and Maunsell,*

574 2009). Neurons in the right hemisphere were considered to be in the attended state when the
 575 attentional cue was on the left, and vice-versa. We note that because our criteria for choosing which
 576 trials and units to analyze were based on different needs for data analysis compared to the original
 577 study (*Cohen and Maunsell, 2009*) the specific firing rates and covariances differ quantitatively
 578 from those previously reported.

579 In monkey 1, an average of 51.1 (min 35, max 80) units were analyzed from the right hemisphere,
 580 and an average of 27.5 (min 14, max 56) units were analyzed from the left hemisphere. From monkey
 581 2, an average of 56.6 (min 43, max 71) units from the right hemisphere, and an average of 37.7 (min
 582 32, max 46) units from the left hemisphere were analyzed. From each recording, spikes falling between
 583 60 and 260 ms from stimulus onset were considered for the firing rate analysis, to account for the
 584 latency of neuronal responses in V4.

585 Comparing change in covariance to change in variance

586 Let S^U be the matrix containing spike counts of the neurons on trials in which they are in the
 587 unattended state, and S^A the matrix containing spike counts of the neurons on trials in which they are
 588 in the attended state. Denote the unattended spike count covariance matrix by $C^U = \text{Cov}(S^U)$, and
 589 the attended one by $C^A = \text{Cov}(S^A)$. Attentional changes in covariance and variance were measured
 590 both on average (Figure 1c) and as distributions (Figure 1d). The distributions of the normalized
 591 differences

$$\frac{\text{Cov}^A - \text{Cov}^U}{\max(|\text{Cov}^A|, |\text{Cov}^U|)} \quad \text{and} \quad \frac{\text{Var}^A - \text{Var}^U}{\max(|\text{Var}^A|, |\text{Var}^U|)} \quad (12)$$

592 reveal a concentration of negative covariance changes, and a distribution of variance changes symmetric
 593 about zero. Here, Cov^A and Cov^U (Var^A and Var^U) are vectors containing covariance (variance) values
 594 of the entire data set. Note that the distributions are bounded between -2 and 2 by construction.

595 Solving systems of equations by error minimization

596 When solving systems of the form of equation (2) in order to quantify the fit of the model, a nonlinear
 597 equation solver (*fminunc*) in MATLAB was used. The solver found minima of an objective function
 598 which we defined as the Euclidean norm of the difference of the approximation of the attended
 599 covariance matrix and the original attended covariance matrix, in other words, the error of the
 600 approximation:

$$f(g_1, \dots, g_N) = \sqrt{\sum_{i < j} (g_i C^U(i, j) g_j - C^A(i, j))^2}. \quad (13)$$

601 Shuffled covariance matrices

602 For finite population sizes ($N < \infty$) we expect our algorithm to extract some low-rank structure
 603 between arbitrary covariance matrices. Let $\sqrt{C^A}$ be the principal square root of the attended covariance
 604 matrix, the unique positive-semidefinite square root of a positive-semidefinite matrix. Consider the
 605 symmetric matrix $D = \text{perm}(\sqrt{C^A})$ computed from a random permutation of the upper-triangular
 606 entries of $\sqrt{C^A}$. Finally, let $C_{\text{shuf}}^A = \text{real}(DD)$. The square root-permutation-squaring procedure
 607 guarantees a positive-semidefinite matrix, as the square of any matrix is positive-semidefinite. Shuffling
 608 removes any relation between C^U and C_{shuf}^A , and any remaining detected structure would be due to
 609 finite sampling. The shuffled covariance gain \hat{g}_{shuf} provides the prediction $\hat{C}_{\text{shuf}}^A := \hat{g}_{\text{shuf}} \hat{g}_{\text{shuf}}^T \circ C^U$,
 610 and ρ_{shuf} measures the relation between \hat{C}_{shuf}^A and C_{shuf}^A . Synthetic data shows that as population
 611 size N becomes large the coefficient ρ_{shuf} approaches 0 (Appendix: Detected structure in random
 612 covariance matrices is a finite-size effect).

613 Upper bound covariance matrices

614 The covariance matrices C^U and C^A are estimates obtained from a finite number of trials, and
 615 any estimation error will compromise the ability to detect rank one structure of A_C . Here we
 616 outline an upper bound for the model performance based on a finite number of trials over which the

617 covariance matrices were originally estimated. Let $\hat{C}^A := \hat{\mathbf{g}}\hat{\mathbf{g}}^T \circ C^U$ with $\hat{\mathbf{g}}$ minimizing the L^2 norm of
 618 $C^A := \mathbf{g}\mathbf{g}^T \circ C^U$. We remark that \hat{C}^A perfectly decomposes according to the statistical model in Eq.
 619 (2). We used \hat{C}^A to generate an artificial set of N correlated Poisson spike counts, using an algorithm
 620 based on a latent multivariate gaussian model (Macke et al., 2009). We sampled these population
 621 spike counts with a fixed number of trials (M) with D be the resulting $M \times N$ matrix of Poisson
 622 samples for each process. Let $C_{ub}^A = \text{Cov}(D)$ be the “upper bound” covariance matrix: a finite trial
 623 sampling approximation to the perfectly decomposable matrix \hat{C}^A . Finally, we employ our algorithm
 624 to give $\hat{C}_{ub}^A := \hat{\mathbf{g}}_{ub}\hat{\mathbf{g}}_{ub}^T C^U$, where the vector $\hat{\mathbf{g}}_{ub}$ minimizes the L^2 norm of the error.

625 Since \hat{C}^A is perfectly decomposable then for $M \rightarrow \infty$ we have $\hat{C}_{ub}^A = C_{ub}^A = \hat{C}^A$. Thus in
 626 the large M limit the coefficient ρ_{ub} between elements of \hat{C}_{ub}^A and C_{ub}^A converges to 1 (Appendix:
 627 Performance limited by available number of trials). However, for finite M we have that $\rho_{ub} < 1$,
 628 solely due to inaccuracies in estimating \hat{C}^A with C_{ub}^A . To account for the possibility of particular
 629 strings of realizations D introducing random biases into C_{ub}^A , we performed the following analysis on
 630 10 independently generated upper-bound covariance matrices C_{ub}^A .

631 Leave-one-out cross-validation

632 Instead of solving the system consisting of all equations (2), we remove one of them. Denote the
 633 complete set of equations by S , an individual equation as $s_{ij} := \{C_{ij}^A = g_i g_j C_{ij}^U\}$ and the set of equations
 634 with one of them removed as $S_{ab} := S - s_{ab}$. We then solve the system S_{ab} . Denote the solution by
 635 \mathbf{g}_{ab} . We can then compare C_{ab}^A and $\hat{C}_{ab}^A = \mathbf{g}_{ab}(a)\mathbf{g}_{ab}(b)C_{ab}^U$. We do this for $\max(1000, N(N-1)/2$
 636 possible systems S_{ab} . The ρ of the vector of resulting C_{ab}^A vs \hat{C}_{ab}^A values is a measure of how well the
 637 system can predict one of its elements, or in other words, how well the structure holds together when
 638 one element is taken out. This leave-one-out cross-validation was performed for the shuffled and the
 639 upper-bound cases as well.

640 Mean field model

641 The mean spiking activity over the population α ($= E$ or I) is

$$r_\alpha(t) = \langle y_{i\alpha}(t) \rangle_i, \quad (14)$$

where $y_{i\alpha}(t) = \sum_{j=1}^{n_{i\alpha}} \delta(t - t_{i\alpha}^j)$ is the spike train of excitatory neuron i of population α , $n_{i\alpha}$ is the
 number of spikes from that neuron, and $t_{i\alpha}^j$ is the time of spike j . We follow previous studies (Tetzlaff
 et al., 2012; Ozeki et al., 2009; Ledoux and Brunel, 2011) and consider the firing rate dynamics
 of the E and I populations given by the system in Eqs. (6):

$$\begin{aligned} \tau_E \frac{dr_E}{dt} &= -r_E + f_E \left(\mu_{EB} + A\Delta\mu_E + J_{EE}r_E - J_{EI}r_I + \sigma_E \left[\sqrt{1-\chi}x_E(t) + \sqrt{\chi}x(t) \right] \right), \\ \tau_I \frac{dr_I}{dt} &= -r_I + f_I \left(\mu_{IB} + A\Delta\mu_I + J_{IE}r_E - J_{II}r_I + \sigma_I \left[\sqrt{1-\chi}x_I(t) + \sqrt{\chi}x(t) \right] \right). \end{aligned}$$

642 Here $\mu_{\alpha\beta}$ is the attention independent drive to population α , $A \in [0, 1]$ is the attention variable,
 643 and $\Delta\mu_\alpha$ is the maximal drive to population α due to attention. The parameter $J_{\alpha\beta}$ is the coupling
 644 from population β to populations α . The stochastic processes $x_E(t)$, $x_I(t)$, and $x(t)$ are the global
 645 fluctuations applied to the network. The excitatory and inhibitory populations have private fluctuations
 646 $x_\alpha(t)$ and also common fluctuations $x(t)$ given to both populations; the parameter χ scales the degree
 647 of private versus common fluctuations. We perform calculations for arbitrary χ and then take $\chi \rightarrow 1$
 648 to match the system given in Eqs. (6). The total intensity of fluctuations to population α is set by σ_α .
 649 These simplified rate equations give an accurate picture of the long-timescale dynamics of networks
 650 of coupled spiking neuron models that are in the fluctuation driven regime (Ledoux and Brunel,
 651 2011). The operative timescale reflects a combination of synaptic and membrane integration; since we
 652 are interested in spiking covariance over time windows that are much longer than these, we take them
 653 to be unity for simplicity.

654 To give a quantitative match between the equilibrium statistics of the rate equations and the leaky
 655 integrate-and-fire (LIF) network simulations we take the transfer function f to be the inverse first

656 passage time of an LIF neuron driven by white noise (*Ledoux and Brunel, 2011*):

$$f_\alpha(I) = \left(\tau_\alpha \sqrt{\pi} \int_{(-V_T+I)/\eta_\alpha}^{(-V_R+I)/\eta_\alpha} \exp(z^2) \operatorname{erfc}(z) dz \right)^{-1}. \quad (15)$$

657 The parameter η_α is the intensity of the external fluctuations given to the LIF neurons (Appendix:
658 Spiking model). The membrane timescale τ gives the dimensions of 1/s to the firing rate r_α . The
659 parameter V_T denotes spike threshold while V_R is the reset potential.

If the input fluctuations, $x(t)$, $x_E(t)$, and $x_I(t)$ are white noise processes then the nonlinearity in f makes the stochastic dynamics of $r_E(t)$ and $r_I(t)$ complicated (non-diffusive). To simplify the analysis we consider $x(t)$ as the limiting process from:

$$\tau_x \frac{dx}{dt} = -x + \sqrt{\tau_x} \xi_x(t),$$

660 for $\tau_x \rightarrow 0$, with $\langle \xi_x(t) \rangle = 0$ and $\langle \xi_x(t) \xi_x(t') \rangle = \delta(t - t')$. This makes $x(t)$ sufficiently smooth in time
661 (the same is true for $x_E(t)$ and $x_I(t)$).

We restrict the coupling $J_{\alpha\beta}$ such that for $\sigma_\alpha = 0$ the equilibrium point (\bar{r}_E, \bar{r}_I) is stable and given by:

$$\begin{aligned} \bar{r}_E &= f_E(\mu_{EB} + A\Delta\mu_E + J_{EE}\bar{r}_E - J_{EI}\bar{r}_I), \\ \bar{r}_I &= f_I(\mu_{IB} + A\Delta\mu_I + J_{IE}\bar{r}_E - J_{II}\bar{r}_I). \end{aligned} \quad (16)$$

For sufficiently small σ_α the fluctuations in population activity about the equilibrium firing rate, $\delta r_\alpha(t) = r_\alpha(t) - \bar{r}_\alpha$, obey the linearized stochastic system:

$$\begin{aligned} \tau_E \frac{d}{dt} \delta r_E &= (-1 + L_E J_{EE}) \delta r_E - L_E J_{EI} \delta r_I + L_E \sigma_E (\sqrt{1-\chi} x_E(t) + \sqrt{\chi} x(t)), \\ \tau_I \frac{d}{dt} \delta r_I &= L_I J_{IE} \delta r_E - (1 + L_I J_{II}) \delta r_I + L_I \sigma_I (\sqrt{1-\chi} x_I(t) + \sqrt{\chi} x(t)). \end{aligned} \quad (17)$$

662 Here $L_\alpha = \frac{df_\alpha}{dI}|_{I=I_\alpha^{\text{eff}}}$ is the slope of the transfer function f_α evaluated at the equilibrium point $I_\alpha^{\text{eff}} =$
663 $\mu_\alpha + A\Delta\mu_\alpha + J_{\alpha E}\bar{r}_E - J_{\alpha I}\bar{r}_I$. Eq. (17) is a two dimensional Ornstein-Uhlenbeck process (*Gardiner, 2004*) that is readily amenable to analysis.

665 Computing V_E

666 In matrix form the system Eq.(17) is written as:

$$\frac{d}{dt} \delta \mathbf{r} = M \delta \mathbf{r} + D \mathbf{x}. \quad (18)$$

667 Here $\delta \mathbf{r} = [\delta r_E, \delta r_I]$, $\mathbf{x} = [x_E, x_I, x]$, and

$$668 M = \begin{bmatrix} -1 + L_E J_{EE} & -L_E J_{EI} \\ L_I J_{IE} & -1 - L_I J_{II} \end{bmatrix} \text{ and } D = \begin{bmatrix} L_E \sigma_E \sqrt{1-\chi} & 0 & L_E \sigma_E \sqrt{\chi} \\ 0 & L_I \sigma_I \sqrt{1-\chi} & L_I \sigma_I \sqrt{\chi} \end{bmatrix}.$$

670 The stationary autocovariance function is computed as:

$$\tilde{C}(s) = \langle \delta \mathbf{r}(t), \delta \mathbf{r}(t+s) \rangle = \begin{cases} \exp(Ms)\Sigma & \text{if } s > 0 \\ \Sigma \exp(-M^T s) & \text{if } s \leq 0 \end{cases}, \quad (19)$$

672 where s is a time lag and $\Sigma = \frac{(\text{Det } M) DD^T + [M - (\text{Tr } M) \mathbf{1}] DD^T [M - (\text{Tr } M) \mathbf{1}]^T}{2(\text{Tr } M)(\text{Det } M)}$ is the variance matrix (Det and
673 Tr denote the determinant and trace operations, respectively). Here, $\mathbf{1}$ is the 2×2 identity matrix.

674 The covariance between populations α and β over long time scales is given by

$$C(\alpha, \beta) = \int_{-\infty}^{\infty} \tilde{C}(s; \alpha, \beta) ds, \quad (20)$$

675 where the integration is performed over the appropriate element of the matrix $\tilde{C}(s)$. In particular, the
676 long timescale variance of the excitatory population is given by (after some algebra):

$$V_E = C(E, E) = \frac{L_E^2}{(1 + J_I L_I - J_E L_E)^2} (J_I L_I (\sigma_E - \sigma_I) + \sigma_E)^2. \quad (21)$$

677 We remark that the long timescale covariance matrix can alternatively be computed from $C =$
678 $M^{-1}D[M^{-1}D]^T$ (*Gardiner, 2004*). To obtain the compact expression for V_E we have assumed
679 symmetric coupling: $J_I := J_{EI} = J_{II}$, $J_E := J_{EE} = J_{IE}$, and $\chi \rightarrow 1$. These are not required for the
680 main results of our study and merely ease the analysis of equations.

681 Computing stimulus response gain

We decompose $\mu_{\alpha B} = k_{\alpha}s + \hat{\mu}_{\alpha B}$ and define the gain of population α to stimulus s as $G_{\alpha} = \frac{d\bar{r}_{\alpha}}{ds} = L_{\alpha} \frac{dI_{\alpha}}{ds}$.
The term $\frac{dI_{\alpha}}{ds}$ is obtained by differentiating equations (16)) with respect to s :

$$\frac{dI_{\alpha}}{ds} = k_{\alpha} + J_E G_E - J_I G_I.$$

682 Solving the system of two equations for G_E yields:

$$G_E = \frac{L_E(k_E + J_I L_I(k_E - k_I))}{1 + J_I L_I - J_E L_E}. \quad (22)$$

683 For the sake of compactness we set $\sigma_E = \sigma_I$ to obtain the result in Eq. (8).

684 Fisher Information

Linear Fisher Information depends on the stimulus response gains and covariance matrix of the excitatory and inhibitory populations:

$$\begin{aligned} \text{FI}_{EI} &= \begin{bmatrix} G_E & G_I \\ C_{EI} & V_I \end{bmatrix} \begin{bmatrix} V_E & C_{EI} \\ C_{EI} & V_I \end{bmatrix}^{-1} \begin{bmatrix} G_E \\ G_I \end{bmatrix} \\ &= \frac{G_E^2 V_I + G_I^2 V_E - 2G_E G_I C_{EI}}{V_E V_I - C_{EI}^2}, \end{aligned} \quad (23)$$

685 When the input correlation $0 \leq \chi < 1$ we have:

$$V_E = \left(\frac{L_E}{1 + J_I L_I - J_E L_E} \right)^2 (J_I^2 L_I^2 (\sigma_E^2 + \sigma_I^2 - 2\sigma_E \sigma_I \chi) + 2J_I L_I \sigma_E (\sigma_E - \sigma_I \chi) + \sigma_E^2), \quad (24)$$

$$V_I = \left(\frac{L_I}{1 + J_I L_I - J_E L_E} \right)^2 (J_E^2 L_E^2 (\sigma_E^2 + \sigma_I^2 - 2\sigma_E \sigma_I \chi) + 2J_E L_E \sigma_I (\sigma_I - \sigma_E \chi) + \sigma_I^2), \quad (25)$$

687 and

$$\begin{aligned} C_{EI} &= \frac{L_E L_I}{(1 + J_I L_I - J_E L_E)^2} (J_E J_I L_E L_I (\sigma_E^2 + \sigma_I^2 - 2\sigma_E \sigma_I \chi) \\ &\quad + J_E L_E \sigma_E (\sigma_E - \sigma_I \chi) - J_I L_I \sigma_I (\sigma_I - \sigma_E \chi) + \sigma_E \sigma_I \chi). \end{aligned} \quad (26)$$

688 Inserting these expressions and those for G_E and G_I into Eq. (23) and simplifying yields:

$$\text{FI}_{EI} = \frac{2\chi k_E k_I \sigma_E \sigma_I - k_E^2 \sigma_I^2 - k_I^2 \sigma_E^2}{(\chi^2 - 1)\sigma_I^2 \sigma_E^2}. \quad (27)$$

689 We remark that FI_{EI} is independent of L_E and L_I and thus independent of attentional modulation.
690 Notice that we have re-introduced the correlation constant χ into the equations, rather than
691 only considering the limit $\chi \rightarrow 1$. If $\chi = 1$, the excitatory and inhibitory populations are receiving
692 completely identical noise. If this is the case, the correlation cancellation would be perfect, leading to
693 infinite informational content, as can be seen in Eq. (27).

694 Acknowledgments

The research was support by National Science Foundation grants NSF-DMS-1313225 (B.D.), NSF DMS-1517082 (B.D.), a grant from the Simons Foundation collaboration on the global brain (SCGB #325293MC;BD), NIH grants 4R00EY020844-03 and R01 EY022930 (MRC), a Whitehall Foundation Grant (MRC), Klingenstein-Simons Fellowship (MRC), a Sloan Research Fellowship (MRC), and a McKnight Scholar Award (MRC). We thank John Maunsell for the generous use of the data, and Kenneth Miller, Ashok Litwin-Kumar, Douglas Ruff, and Robert Rosenbaum for useful discussions.

Table 1. Model Parameters

Parameter	Description	Value
τ	Time constants for membrane dynamics	0.01 s
V_T	Spike Threshold	1
V_R	Spike Reset	0
μ_E	Excitatory baseline bias	0.6089
μ_I	Inhibitory baseline bias	0.5388
$\Delta\mu_E$	Attentional modulation of excitatory bias	0.2624
$\Delta\mu_I$	Attentional modulation of inhibitory bias	0.3608
J_E	Excitatory coupling constant	1.5
J_I	Inhibitory coupling constant	3
σ_E	Amplitude of external noise to E population	0.3
σ_I	Amplitude of external noise to I population	0.35
c	Proportion of common noise to E and I populations	1
k_E	Sensitivity of E population to stimulus input	1
k_I	Sensitivity of I population to stimulus input	0

References

- 701 **References**
- 702 Alitto HJ, Dan Y. Cell-type-specific modulation of neocortical activity by basal forebrain input. *Frontiers in*
 703 *systems neuroscience*. 2012; 6.
- 704 Ardid S, Wang XJ, Compte A. An integrated microcircuit model of attentional processing in the neocortex.
 705 *Journal of Neuroscience*. 2007; 27(32):8486–8495.
- 706 Averbeck BB, Latham PE, Pouget A. Neural correlations, population coding and computation. *Nature*
 707 *reviews neuroscience*. 2006; 7:358–366.
- 708 Beck J, Bejjani VR, Pouget A. Insights from a simple expression for linear fisher information in a recurrently
 709 connected population of spiking neurons. *Neural computation*. 2011; 23(6):1484–1502.
- 710 Ben-Yishai R, Bar-Or RL, Sompolinsky H. Theory of orientation tuning in visual cortex. *Proceedings of the*
 711 *National Academy of Sciences*. 1995; 92(9):3844–3848.
- 712 Buia CI, Tiesinga PH. Role of interneuron diversity in the cortical microcircuit for attention. *Journal of*
 713 *neurophysiology*. 2008; 99(5):2158–2182.
- 714 Cardin JA, Palmer LA, Contreras D. Stimulus Feature Selectivity in Excitatory and Inhibitory Neurons in
 715 Primary Visual Cortex. *Journal of Neuroscience*. 2007; 27(39):10333–10344.
- 716 Cohen MR, Maunsell JHR. Attention improves performance primarily by reducing interneuronal correlations.
 717 *Nature neuroscience*. 2009; 12(12):1594–1600.
- 718 Cohen MR, Maunsell JHR. Using neuronal populations to study the mechanisms underlying spatial and
 719 feature attention. *Neuron*. 2011; 70(6):1192–1204.
- 720 Crochet S, Poulet JF, Kremer Y, Petersen CC. Synaptic mechanisms underlying sparse coding of active touch.
 721 *Neuron*. 2011; 69(6):1160–1175.
- 722 Deco G, Thiele A. Cholinergic control of cortical network interactions enables feedback-mediated attentional
 723 modulation. *European Journal of Neuroscience*. 2011; 34(1):146–157.
- 724 Disney AA, Alasady HA, Reynolds JH. Muscarinic acetylcholine receptors are expressed by most parvalbumin-
 725 immunoreactive neurons in area MT of the macaque. *Brain and Behavior*. 2014; 4(3):431–445.
- 726 Disney AA, Aoki C, Hawken MJ. Gain modulation by nicotine in macaque V1. *Neuron*. 2007; 56:701–713.
- 727 Doiron B, Lindner B, Longtin A, Maler L, Bastian J. Oscillatory activity in electrosensory neurons increases
 728 with the spatial correlation of the stochastic input stimulus. *Physical Review Letters*. 2004; 93(4).

- 729 Doiron B, Litwin-Kumar A, Rosenbaum R, Ocker G, Josic K. The mechanics of state dependent neural
730 correlations. *Nature Neuroscience*. 2016; 19(1):383–393.
- 731 Ecker AS, Denfield GH, Bethge M, Tolias AS. On the Structure of Neuronal Population Activity under
732 Fluctuations in Attentional State. *Journal of Neuroscience*. 2016 Feb; 36(5):1775–1789.
- 733 Ecker AS, Berens P, Cotton RJ, Subramaniyan M, Denfield GH, Cadwell CR, Smirnakis SM, Bethge M,
734 Tolias AS. State dependence of noise correlations in macaque primary visual cortex. *Neuron*. 2014 Apr;
735 82(1):235–248.
- 736 Ecker AS, Denfield GH, Bethge M, Tolias AS. On the structure of population activity under fluctuations in
737 attentional state. *bioRxiv*. 2015; p. 018226.
- 738 Engel TA, Steinmetz NA, Gieselmann MA, Thiele A, Moore T, Boahen K. Selective modulation of cortical
739 state during spatial attention. *Science*. 2016; 354(6316):1140–1144.
- 740 Fries P, Reynolds JH, Rorie AE, Desimone R. Modulation of Oscillatory Neuronal Synchronization by Selective
741 Visual Attention. *Science*. 2001; 291:1560–1563.
- 742 Fu Y, Tucciarone JM, Espinosa JS, Sheng N, Darcy DP, Nicoll RA, Huang ZJ, Stryker MP. A cortical circuit
743 for gain control by behavioral state. *Cell*. 2014; 156(6):1139–1152.
- 744 Gardiner CW. *Handbook of Stochastic Methods for Physics, Chemistry and the Natural Sciences*. 3rd ed.
745 Springer-Verlag; 2004.
- 746 Gilbert CD, Sigman M. Brain states: top-down influences in sensory processing. *Neuron*. 2007; 54(5):677–696.
- 747 Ginzburg I, Sompolinsky H. Theory of correlations in stochastic neural networks. *Physical Review E*. 1994;
748 50(4):3171–3191.
- 749 Harris KD, Thiele A. Cortical state and attention. *Nature reviews neuroscience*. 2011; 12(9):509–523.
- 750 Hasselmo ME. Neuromodulation and cortical function: modeling the physiological basis of behavior.
751 Behavioural brain research. 1995; 67(1):1–27.
- 752 Helias M, Tetzlaff T, Diesmann M. The correlation structure of local neuronal networks intrinsically results
753 from recurrent dynamics. *PLoS Comput Biol*. 2014; 10(1).
- 754 Herrero JL, Gieselmann MA, Sanayei M, Thiele A. Attention-induced variance and noise correlation reduction
755 in macaque V1 is mediated by NMDA receptors. *Neuron*. 2013; 78(4):729–739.
- 756 Isaacson JS, Scanziani M. How inhibition shapes cortical activity. *Neuron*. 2011; 72(2):231–243.
- 757 Kelly RC, Smith MA, Kass RE, Lee TS. Local field potentials indicate network state and account for neuronal
758 response variability. *Journal of computational neuroscience*. 2010; 29(3):567–579.
- 759 Kepecs A, Fishell G. Interneuron cell types are fit to function. *Nature*. 2014; 505(7483):318–326.
- 760 Kim H, Ährlund-Richter S, Wang X, Deisseroth K, Carlén M. Prefrontal Parvalbumin Neurons in Control of
761 Attention. *Cell*. 2016; 164(1):208–218.
- 762 Kohn A, Coen-Cagli R, Kanitscheider I, Pouget A. Correlations and neuronal population information. *Annual
763 review of neuroscience*. 2016; 39:237–256.
- 764 Kuchibhotla KV, Gill JV, Lindsay GW, Papadoyannis ES, Field RE, Sten TAH, Miller KD, Froemke RC.
765 Parallel processing by cortical inhibition enables context-dependent behavior. *Nature Neuroscience*. 2017;
766 20(1):62–71.
- 767 Ledoux E, Brunel N. Dynamics of networks of excitatory and inhibitory neurons in response to time-dependent
768 inputs. *Frontiers in Computational Neuroscience*. 2011; 5.
- 769 Lee SH, Dan Y. Neuromodulation of brain states. *Neuron*. 2012; 76(1):209–222.
- 770 Lee S, Kruglikov I, Huang ZJ, Fishell G, Rudy B. A disinhibitory circuit mediates motor integration in the
771 somatosensory cortex. *Nature neuroscience*. 2013; 16(11):1662–1670.
- 772 Lin IC, Okun M, Carandini M, Harris KD. The Nature of Shared Cortical Variability. *Neuron*. 2015;
773 87(3):644–656.

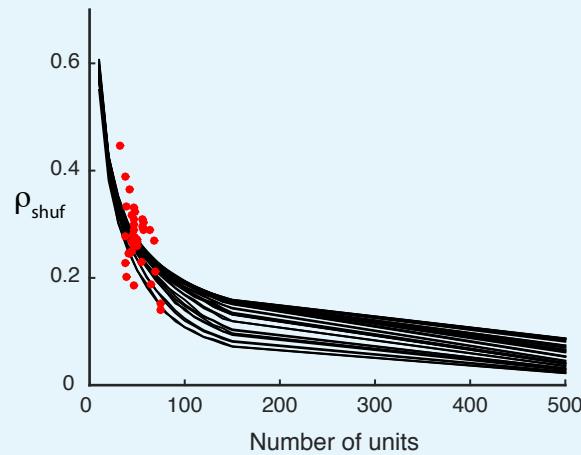
- 774 Ly C, Middleton JW, Doiron B. Cellular and circuit mechanisms maintain low spike co-variability and enhance
775 population coding in somatosensory cortex. *Frontiers in Computational Neuroscience*. 2012; 6.
- 776 Macke JH, Berens P, Ecker AS, Tolias AS, Bethge M. Generating spike trains with specified correlation
777 coefficients. *Neural Computation*. 2009 February; 21(2):397–423.
- 778 McAdams CJ, Maunsell JH. Attention to both space and feature modulates neuronal responses in macaque
779 area V4. *Journal of Neurophysiology*. 2000; 83(3):1751–1755.
- 780 Mitchell JF, Sundberg KA, Reynolds JH. Differential attention-dependent response modulation across cell
781 classes in macaque visual area V4. *Neuron*. 2007; 55(1):131–141.
- 782 Mitchell JF, Sundberg KA, Reynolds JH. Spatial attention decorrelates intrinsic activity fluctuations in
783 macaque area V4. *Neuron*. 2009; 63(6):879–888.
- 784 Mongillo G, Hansel D, van Vreeswijk C. Bistability and spatiotemporal irregularity in neuronal networks
785 with nonlinear synaptic transmission. *Physical review letters*. 2012; 108(15):158101.
- 786 Moore T, Zirnsak M. Neural Mechanisms of Selective Visual Attention. *Annual Review of Psychology*. 2017;
787 68:47–72.
- 788 Moreno-Bote R, Beck J, Kanitscheider I, Pitkow X, Latham P, Pouget A. Information-limiting correlations.
789 *Nature neuroscience*. 2014; 17(10):1410–1417.
- 790 Navalpakkam V, Itti L. Modeling the influence of task on attention. *Vision research*. 2005; 45(2):205–231.
- 791 Noudoost B, Moore T. The role of neuromodulators in selective attention. *Trends in cognitive sciences*. 2011;
792 15(12):585–591.
- 793 Ocker GK, Hu Y, Buice MA, Doiron B, Josić K, Rosenbaum R, Shea-Brown E. From the statistics of
794 connectivity to the statistics of spike times in neuronal networks. *arXiv preprint arXiv:170303132*. 2017; .
- 795 Ozeki H, Finn IM, Schaffer ES, Miller KD, Ferster D. Inhibitory stabilization of the cortical network underlies
796 visual surround suppression. *Neuron*. 2009; 62(4):578–592.
- 797 Pernice V, Staude B, Cardanobile S, Rotter S. How structure determines correlations in neuronal networks.
798 *PLoS Comput Biol*. 2011; 7(5).
- 799 Pfeffer CK, Xue M, He M, Huang ZJ, Scanziani M. Inhibition of inhibition in visual cortex: the logic of
800 connections between molecularly distinct interneurons. *Nature neuroscience*. 2013; 16(8):1068–1076.
- 801 Poorthuis RB, Enke L, Letzkus JJ. Cholinergic circuit modulation through differential recruitment of
802 neocortical interneuron types during behaviour. *The Journal of physiology*. 2014; 592(19):4155–4164.
- 803 Pouget A, Dayan P, Zemel R. Information processing with population codes. *Nature Reviews Neuroscience*.
804 2000; 1(2):125–132.
- 805 Rabinowitz NC, Goris RL, Cohen M, Simoncelli E. Attention stabilizes the shared gain of V4 populations.
806 *eLife*. 2015; p. e08998.
- 807 Renart A, De la Rocha J, Bartho P, Hollender L, Parga N, Reyes A, Harris KD. The asynchronous state in
808 cortical circuits. *Science*. 2010; 327(5965):587–590.
- 809 Reynolds JH, Chelazzi L. Attentional modulation of visual processing. *Annu Rev Neurosci*. 2004; 27:611–647.
- 810 Reynolds JH, Chelazzi L, Desimone R. Competitive mechanisms subserve attention in macaque areas V2
811 and V4. *The Journal of Neuroscience*. 1999; 19(5):1736–1753.
- 812 Reynolds JH, Heeger DJ. The normalization model of attention. *Neuron*. 2009; 61(2):168–185.
- 813 Rosenbaum R, Smith MA, Kohn A, Rubin JE, Doiron B. The spatial structure of correlated neuronal
814 variability. *Nature Neuroscience*. 2017; 20(1):107–114.
- 815 Rubin DB, Van Hooser SD, Miller KD. The stabilized supralinear network: a unifying circuit motif underlying
816 multi-input integration in sensory cortex. *Neuron*. 2015; 85(2):402–417.
- 817 Rudy B, Fishell G, Lee S, Hjerling-Leffler J. Three groups of interneurons account for nearly 100% of
818 neocortical GABAergic neurons. *Developmental neurobiology*. 2011; 71(1):45–61.

- 819 Ruff DA, Cohen MR. Attention can either increase or decrease spike count correlations in visual cortex.
820 Nature neuroscience. 2014; 17(11):1591–1597.
- 821 Sanayei M, Herrero J, Distler C, Thiele A. Attention and normalization circuits in macaque V1. European
822 Journal of Neuroscience. 2015; 41(7):949–964.
- 823 Silver RA. Neuronal arithmetic. Nature Reviews Neuroscience. 2010; 11(7):474–489.
- 824 Sohya K, Kameyama K, Yanagawa Y, Obata K, Tsumoto T. GABAergic neurons are less selective to stimulus
825 orientation than excitatory neurons in layer II/III of visual cortex, as revealed by *in vivo* functional Ca²⁺
826 imaging in transgenic mice. The Journal of neuroscience. 2007; 27(8):2145–2149.
- 827 Steriade M, McCormick DA, Sejnowski TJ. Thalamocortical oscillations in the sleeping and aroused brain.
828 Science. 1993; 262(5134):679–685.
- 829 Stringer C, Pachitariu M, Steinmetz NA, Okun M, Bartho P, Harris KD, Sahani M, Lesica NA. Inhibitory
830 control of correlated intrinsic variability in cortical networks. Elife. 2016; 5:e19695.
- 831 Sutherland C, Doiron B, Longtin A. Feedback-induced gain control in stochastic spiking networks. Biological
832 cybernetics. 2009; 100(6):475–489.
- 833 Tetzlaff T, Helias M, Einevoll GT, Diesmann M. Decorrelation of Neural-Network Activity by Inhibitory
834 Feedback. PLoS Comput Biol. 2012; 8(8).
- 835 Treue S. Neural correlates of attention in primate visual cortex. Trends in neurosciences. 2001; 24(5):295–300.
- 836 Trousdale J, Hu Y, Shea-Brown E, Josić K. Impact of Network Structure and Cellular Response on Spike
837 Time Correlations. PLoS Comput Biol. 2012; 8(3).
- 838 van Vreeswijk C, Sompolinsky H. Chaotic balanced state in a model of cortical circuits. Neural computation.
839 1998; 10(6):1321–1371.
- 840 Whiteway MR, Butts DA. Revealing unobserved factors underlying cortical activity with a rectified latent
841 variable model applied to neural population recordings. Journal of neurophysiology. 2017; 117(3):919–936.
- 842 Williford T, Maunsell JH. Effects of spatial attention on contrast response functions in macaque area V4.
843 Journal of Neurophysiology. 2006; 96(1):40–54.

844 Appendix 1

845
 846
 847
 848
 849
 850
Detected structure in random covariance matrices is a finite-size effect
 851 Here we show that any prediction of rank one structure in our shuffled covariance matrix
 852 (non-zero ρ_{shuf} in Fig. 2 of the main text) is a finite-data effect. The trial-by-trial covariance
 853 matrices of the experimental data are computed from the spike counts recorded from a set
 854 number of units. To explore the effect of population size on the detected structure in the
 855 shuffled covariance matrices we must rely on synthetic data.
 856

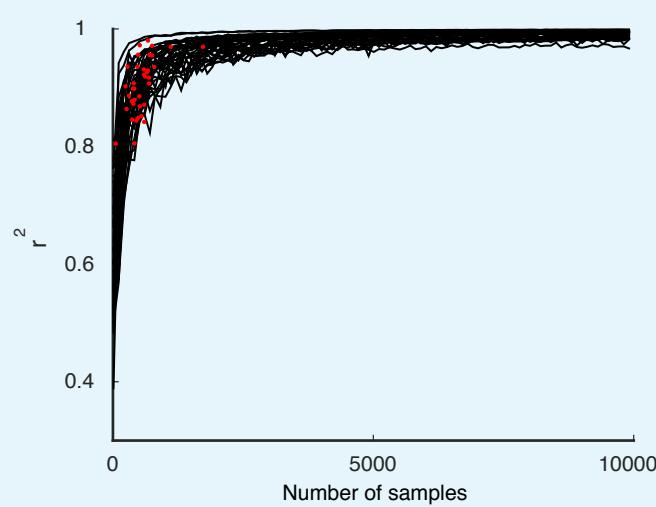
We construct the synthetic covariance matrices by generating Gaussian random numbers with the same mean and standard deviation as the actual covariance matrices from the data. This construction serves as a substitute for the shuffled covariance matrices, and allows for arbitrarily large populations. As we increase the number of units from near 10 to 500, ρ_{shuf} decreases accordingly, indicating that any positive ρ_{shuf} is due to the finite population size, rather than any inherent structure in the data (Appendix Fig. 1).



857
 858
 859
 860
 861
Appendix 1 Figure 1. Detected structure in randomly generated covariance matrices is a finite-size
 862 effect. The model performance (ρ_{shuf}) decreases with increasing system size (black curves). The ρ_{shuf}
 computed from the shuffled neural data (red dots) falls in the same area as the synthetic data
 performance, suggesting that the synthetic data is a reasonable stand-in.

863 **Model performance is limited by number of trials in data**

864 The upper bound for our model ρ_{ub} did not saturate 1 (see Fig. 2 of the main text). Here,
 865 we show that this is also due the finite data available. If infinitely many trials were available
 866 to compute the spike count covariance matrices from the data, and the data obeyed by the
 867 low-rank statistical model, the performance of the model (ρ_{ub}) should tend to one. To test
 868 this, we generate synthetic data from correlated Poisson processes as in the upper bound
 869 computation of the main text but do not limit the number of samples to the number of trials in
 870 the original data. As the number of samples increases we find that $\rho_{\text{ub}} \rightarrow 1$ (Appendix Fig. 2).



871
872
873
875

Appendix 1 Figure 2. The performance of the model ρ_{ub} (black curves) on synthetic data using increasing numbers of Poisson realizations approaches 1. The Poisson model computed with the same number of trials as the data is shown for comparison (red dots).

876
877
878

Model performance for all monkeys and hemispheres

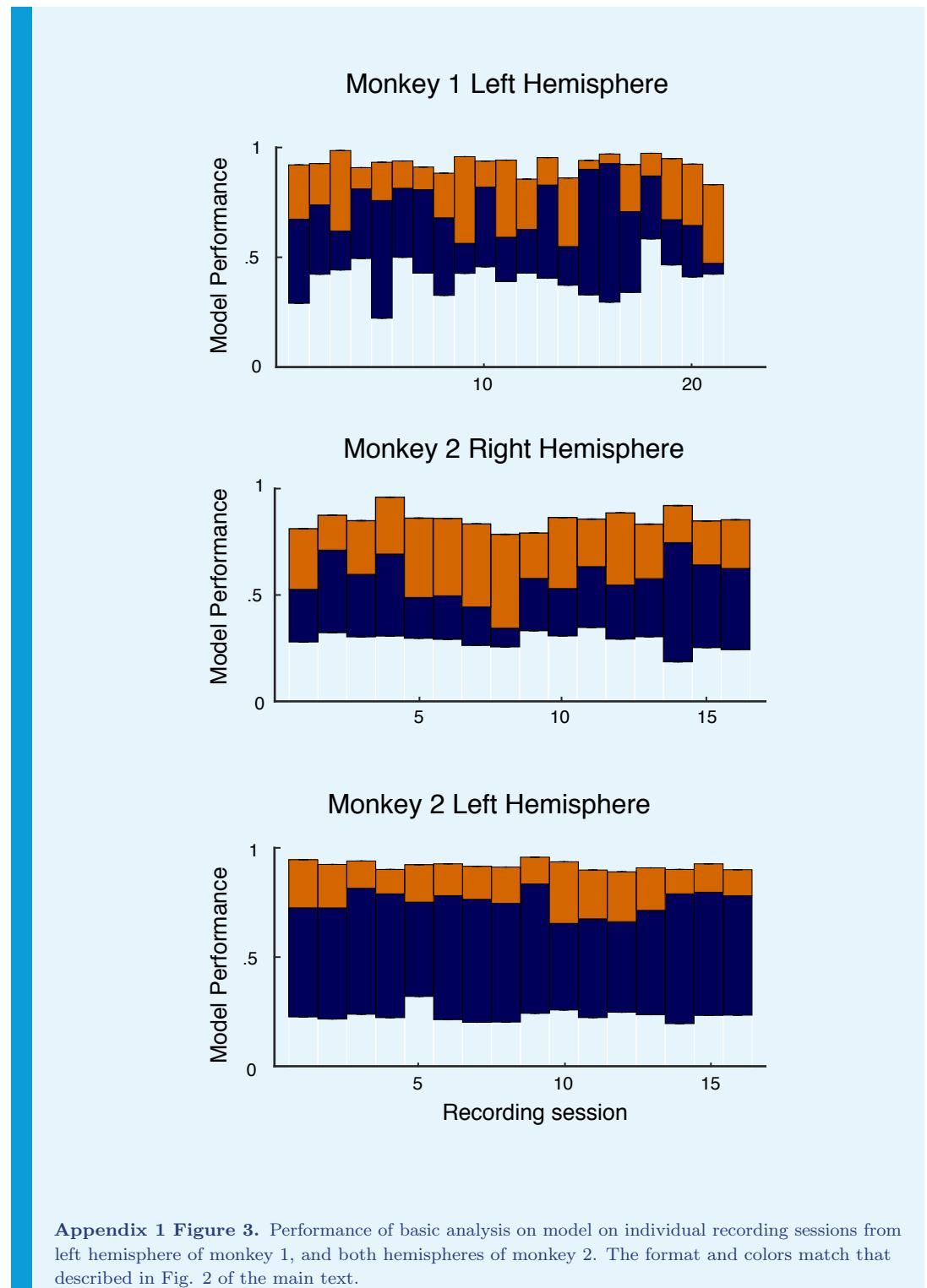
The model performance for individual recording sessions are given here for transparency (Appendix Fig. 3 for the full data and Appendix 4 for the leave-one-out cross validation).

879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898

Low-dimensional modulation is intrinsic to neurons

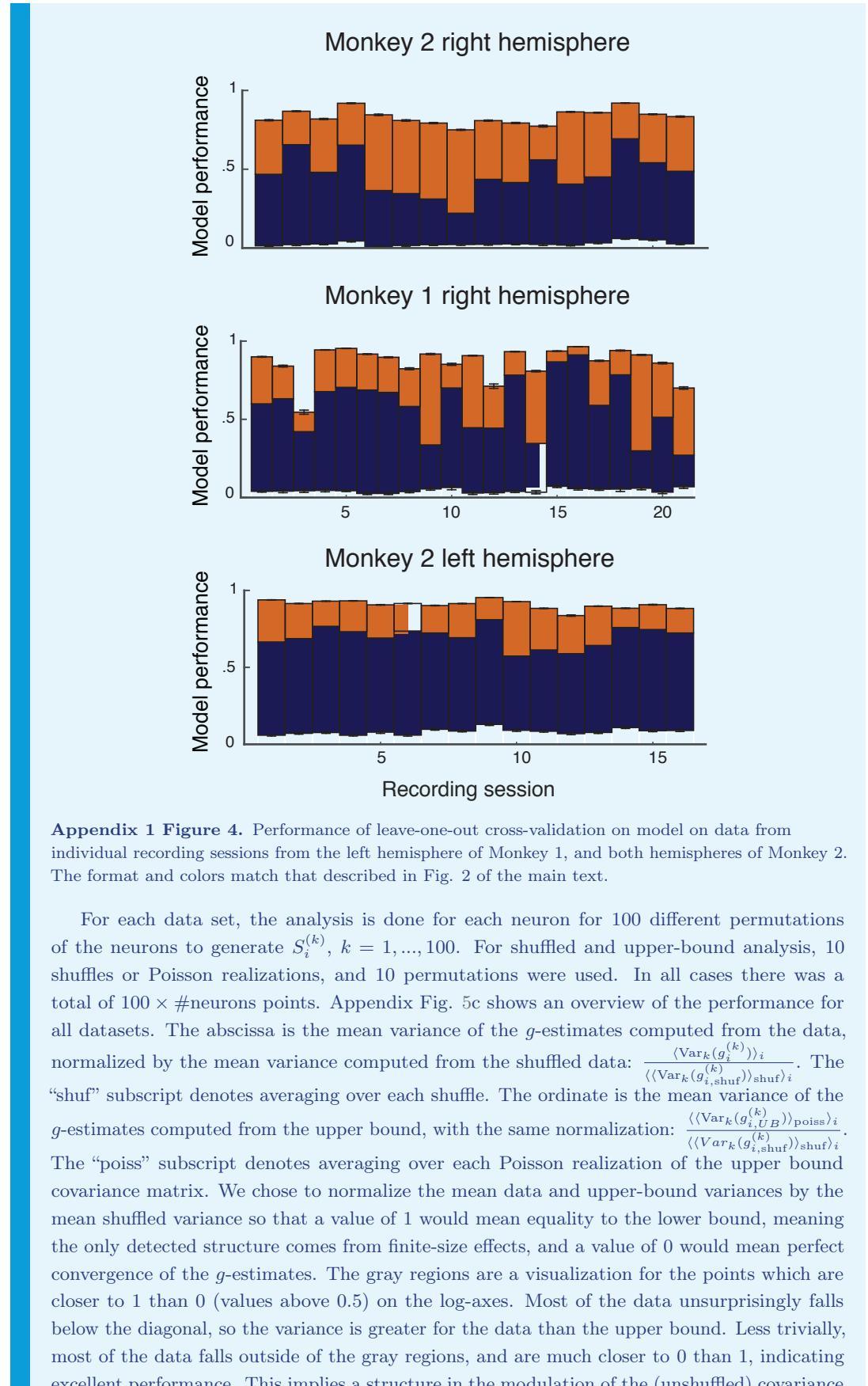
In order to further test our model, we asked to what extent the actual value of the covariance gain g_i of neuron i depends on the neural population whose covariance matrix g_i was estimated from. If we had solved the system S of equations $C_{i,j}^A = g_i g_j C_{i,j}^U$ using covariance matrices computed from recordings from a different set of neurons (including neuron i), would the value of g_i be different? If not, this would be further indication of the independence of the attentional modulation of neuron i from the particular set of other neurons it is analyzed with.

We tackle this question by dividing a set of N neurons into k sets $S_i^{(1)}, S_i^{(2)}, \dots, S_i^{(k)}$ of $m \equiv (N+1)/2$ neurons each that all contain the neuron n_i ($m \equiv N/2+1$ if N is originally even). As an example take $k=2$ and consider the set of neurons n_1, \dots, n_{2i-1} partitioned into two subsets $S_i^{(1)} = \{n_1, \dots, n_i\}$ and $S_i^{(2)} = \{n_i, \dots, n_{2i-1}\}$ (Appendix Fig. 5a). We solve Eq. (1) using the systems of equations obtained from $S_i^{(1)}$ and $S_i^{(2)}$, and obtain two solutions $\mathbf{g}_i^{(1)}$ and $\mathbf{g}_i^{(2)}$. We take the variance of the g -estimations as a metric for how closely the different subsets can estimate an intrinsic value of g . A higher variance would indicate a poorer convergence, and therefore a lower degree of independence from other neurons. Appendix Fig. 5b shows the spread of g -estimates from one dataset for the data, as well as the upper (UB) and lower (shuf) bounds. This spread includes estimates for all g -values for all neurons. The spread in the shuffled case (SEM= 7.42) is largest by two orders of magnitude, and the spread of the upper bound (SEM= 2.60×10^{-3}) is only one order of magnitude tighter than that of the data (SEM= 1.03×10^{-2}), so this case is close to ideal.

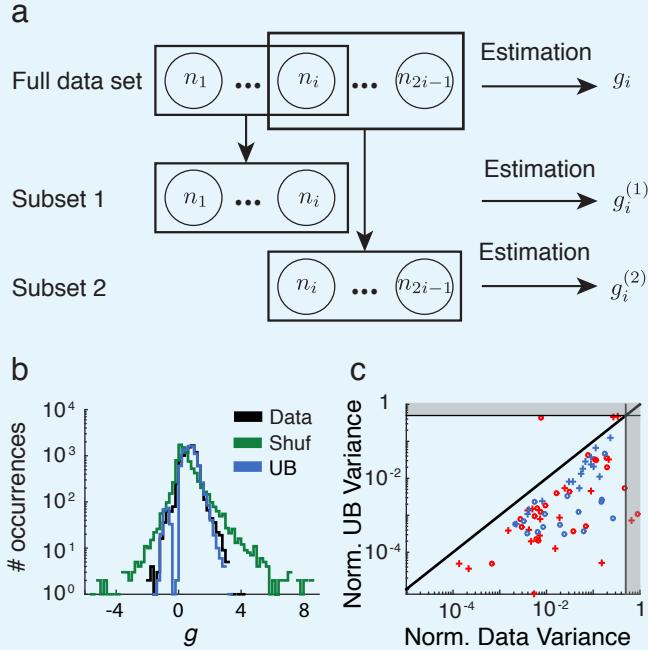


899
900
901
903

Appendix 1 Figure 3. Performance of basic analysis on model on individual recording sessions from left hemisphere of monkey 1, and both hemispheres of monkey 2. The format and colors match that described in Fig. 2 of the main text.



matrices that is preserved over analysis in the contexts of different groups of other neurons. In other words, attention modulates the individual neurons to a large extent independently, in a low-dimensional manner.



Appendix 1 Figure 5. Overlap analysis of gain parameters. **a**, Schematic of overlap analysis. A set of n_{2i-1} neurons is divided into two sets S_1 and S_2 of i , which overlap by exactly one neuron, indexed without loss of generality as neuron i . Parameter g_i is computed using S_1 and S_2 , resulting in two estimates $g_i^{(1)}$ and $g_i^{(2)}$. **b**, Spread of g estimates for the data (black), as well as the upper (blue) and lower (green) bounds, from one day of recordings in one monkey. **c**, Mean variance of the g estimates computed from the data (abscissa) vs from the upper bound (ordinate), normalized by the mean shuffled variance. Each color denotes one of the monkeys, circles denote the right hemisphere recordings, and pluses denote the left hemisphere recordings. The gray regions consist of those points that are beyond 0.5, and therefore closer to the lower bound than the upper bound.

929
930
931
932
933
934
935
936
937
938
939

Network requirements for attentional modulation

In this section we study a network of N neurons with the spike train output from neuron i being $y_i(t) = \sum_k \delta(t - t_{ik})$ where t_{ik} is the k^{th} spike time from neuron i . We consider multiple trials of the discrimination experiment and model the spike train only over a time period $t \in (0, T)$, where we assume that the spike trains to have reached equilibrium statistics. We abuse notation and take the spike count from neuron i over a trial as $y_i = \int_0^T y_i(t) dt$. The trial-to-trial covariance matrix of the network response is \mathbf{C} with element $c_{ij} = \text{Cov}(y_i, y_j)$.

To analyze the network activity we first assume that each spike train is simply perturbed about a background state and employ the linear response ansatz (*Ginzburg and Sompolinsky, 1994; Doiron et al., 2004; Trousdale et al., 2012*) :

$$y_i = y_{iB} + L_i \left(\sum_{k=1}^N J_{ik} y_k + \xi_i \right). \quad (28)$$

Here, J_{ik} is the synaptic coupling from neuron k to neuron i (proportional to the synaptic weight), and ξ_i is a fluctuating external input given to neuron i . The background state of neuron i is y_{iB} , and it represents the stochastic output of a neuron that is not due to the recurrence from the network ($J = 0$) or the external input ($\xi_i = 0$). Finally, L_i is the input to output gain of a neuron i . In this framework, y_i , y_{iB} , and ξ_i are random variables, while L_i and J_{ik} are parameters that describe the intrinsic and network properties of the system.

958
959
960
961
962

Without loss of generality we take $\langle y_{iB} \rangle = 0$, $\langle \xi_i \rangle = 0$, making $\langle y_i \rangle = 0$ a solution for the mean activity. We remark that formally Eq. (28) is incorrect as written; y_i is a random integer while, for instance, $L_i J_{ik} y_k$ need not be an integer. Eq. (28) is only correct upon taking an expectation (over trials) of y_i .

963
964

Here we derive the requirements for external fluctuations and internal coupling for network covariability \mathbf{C} to satisfy the following two conditions (on average):

965
966

C1: $c_{ij}^A = g_i g_j c_{ij}^U$; attentional modulation of covariance is rank one.

C2: $g_i < 1$; spike count covariance decreases with attention.

967
968
969

It is convenient to write Eq. (28) in matrix form and isolate for the population response:

$$\vec{y} = (\mathbf{I} - \mathbf{K})^{-1} (\vec{y}_B + \mathbf{L} \vec{\xi}). \quad (29)$$

970
971
972
973

Here $\vec{y} = [y_1, \dots, y_N]^T$ with similar notation for \vec{y}_B and $\vec{\xi}$. The matrix \mathbf{K} has element $\mathbf{K}_{ij} = L_i J_{ij}$, while $\mathbf{L} = \text{diag}(L_i)$ and \mathbf{I} is the identity matrix. Using Eq. (29) we can express the covariance matrix $\mathbf{C} = \langle \vec{y} \vec{y}^T \rangle$ as:

974
975
976

$$\mathbf{C} = \underbrace{(\mathbf{I} - \mathbf{K})^{-1} \mathbf{B} (\mathbf{I} - \mathbf{K}^T)^{-1}}_{\text{internal covariability}} + \underbrace{(\mathbf{I} - \mathbf{K})^{-1} \mathbf{L} \mathbf{X} \mathbf{L}^T (\mathbf{I} - \mathbf{K}^T)^{-1}}_{\text{external covariability}}, \quad (30)$$

977
978
979
980

where T denotes the transpose operation. Here $\mathbf{B} = \langle \vec{y}_B \vec{y}_B^T \rangle$ is the background covariance, which we take to be simply $\mathbf{B} = \text{diag}(b_i)$. The input covariance matrix is $\mathbf{X} = \langle \vec{\xi} \vec{\xi}^T \rangle$ with elements x_{ij} . In the above we assumed that $\langle \vec{y}_B \vec{\xi}^T \rangle = \mathbf{0}$, meaning that the background state is uncorrelated with the external noisy input.

981
982
983
984

It is clear that \mathbf{C} naturally decomposes into two terms. The first term represents the correlations that are internally generated within the network, via the direct synaptic coupling \mathbf{K} acting upon the background state \mathbf{B} . The second term is how the direct synaptic coupling \mathbf{K} filters the externally applied correlations \mathbf{X} .

986

Satisfying C1

988
989
990
991
992

The background matrix \mathbf{B} is a diagonal matrix and is hence rank N . The high rank \mathbf{B} combined with attentional modulations of both \mathbf{B} and \mathbf{K} make it impossible to satisfy condition **C1**. If the spectral radius of \mathbf{K} is less than 1, then we can expand $(\mathbf{I} - \mathbf{K})^{-1} = \mathbf{I} + \sum_{n=1}^{\infty} \mathbf{K}^n$ (*Pernice et al., 2011; Trousdale et al., 2012*). Inserting this expansion into the expression for the internally generated covariability yields:

993
994

$$(\mathbf{I} - \mathbf{K})^{-1} \mathbf{B} (\mathbf{I} - \mathbf{K}^T)^{-1} = \mathbf{B} + \mathbf{B} \mathbf{K}^T + \mathbf{K} \mathbf{B} + \mathbf{K} \mathbf{B} \mathbf{K}^T + \dots.$$

995
996

Extracting the covariance between neuron i and j ($i \neq j$) due to internal coupling within the network gives:

997
998

$$c_{ijB} = b_i L_j J_{ji} + b_j L_i J_{ij} + \sum_k L_i L_j b_k J_{ik} J_{jk} + \dots.$$

999
1000
1001
1002

If we take $J_{ij} \sim 1/N$ and the network connectivity to be dense (meaning the connection probability is $\sim \mathcal{O}(1)$) then each term is $\mathcal{O}(1/N)$. So long as the spectral radius of \mathbf{K} is less than 1 then the series converges and as $N \rightarrow \infty$ we have that c_{ijB} vanishes (*Pernice et al., 2011; Trousdale et al., 2012; Helias et al., 2014*).

1003
1004
1005
1006
1007
1008

This argument can be extended to networks with $J_{ij} \sim 1/\sqrt{N}$ when combined with a balance condition between recurrent excitation and inhibition. Such networks also produce an asynchronous state where $c_{ijB} \sim 1/N$, vanishing in the large N limit (*Renart et al., 2010*). However, formally balanced networks in the asynchronous state with $N \rightarrow \infty$ have solutions that do not depend on the firing rate transfer L . The attention dependent modulation $A_L : L^U \rightarrow L^A$ is a critical component of our model and care must be taken in ensuring that

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

In contrast, the external covariance \mathbf{X} is not a diagonal matrix, so that the contributions from external fluctuations to \mathbf{C} scale as $N^2 J^2$. This is $\mathcal{O}(1)$ for $J \propto 1/N$. Thus, while the terms in \mathbf{X} must be weak for the linear approximation in Eq. (28) to hold, they need not vanish for large N . Indeed, for moderate \mathbf{X} and large network size it is reasonable to ignore the contribution of internally generated fluctuations to \mathbf{C} . Recent analysis of cortical population recordings show that the shared spiking variability across the population can be well approximated by a rank one model of covariability (Ecker et al., 2014; Lin et al., 2015; Ecker et al., 2015; Rabinowitz et al., 2015). Thus motivated, we take the external fluctuations $\mathbf{X} = \mathbf{x}\mathbf{x}^T$ where $\mathbf{x} = [x_1, \dots, x_N]^T$. In total, we have for large N the approximation:

$$\mathbf{C} \approx ((\mathbf{I} - \mathbf{K})^{-1} \mathbf{Lx}) ((\mathbf{I} - \mathbf{K})^{-1} \mathbf{Lx})^T = \mathbf{cc}^T. \quad (31)$$

Hence \mathbf{C} is rank one matrix with $\mathbf{c} = ((\mathbf{I} - \mathbf{K})^{-1} \mathbf{Lx}) = [c_1, \dots, c_N]^T$. It is trivial to satisfy condition **C1** with $g_i = c_i^A / c_i^U$.

Satisfying C2

We again use the expansion $(\mathbf{I} - \mathbf{K})^{-1} = \mathbf{I} + \sum_{n=1}^{\infty} \mathbf{K}^n$. Truncating this expansion at $n = 1$ yields an approximation considering only synaptic paths of length 1 in the network, and neglecting higher order paths. This is appropriate for J_{ij} sufficiently small. Truncating after inserting the expansion into Eq. (30) yields the following approximation for \mathbf{c} :

$$\mathbf{c} \approx (\mathbf{I} + \mathbf{K}) \mathbf{Lx}. \quad (32)$$

The analysis in the main text begins with this approximation to derive Eq. (5) of the main text.

Spiking Network

Spiking Network Description

We implement a network of leaky integrate-and-fire neurons (LIF) with 1000 excitatory neurons and 200 inhibitory neurons. Individual neurons were modeled as integrate-and-fire units whose voltages obeyed

$$\frac{dV_i}{dt} = \frac{1}{\tau} (\mu_i - V_i) + I_i^{\text{syn}} + I_i^{\text{ext}} \quad (33)$$

for neuron i . When the voltage reached a threshold $V_{\text{th}} = 1$, a spike was recorded and the voltage reset to $V_{\text{re}} = 0$. Time was measured in units of the membrane time constant, $\tau = 1$ for all neurons. The bias μ depended on neuron type and attentional state. In the unattended state, the bias for excitatory neurons was $\mu_E^{\text{un}} = .6089$ and $\mu_I^{\text{un}} = .5388$. In the attended state, $\mu_E^{\text{att}} = .8713$ and $\mu_I^{\text{att}} = .8996$. The recurrent input to neuron i was

$$I_i^{\text{syn}}(t) = \sum_j \mathbf{W}_{ij} \mathbf{J}_{ij}(t) * y_j(t) \quad (34)$$

where \mathbf{W}_{ij} is the strength of the connection from neuron j to neuron i , $J_{ij}(t)$ is the synaptic filter for the projection from neuron j to neuron i , $*$ denotes convolution and $y_j(t)$ is neuron j 's spike train – a series of δ -functions centered at spike times. The synaptic filters were taken to be alpha functions,

$$J_{ij}(t) = \frac{t}{\tau_s} e^{-t/\tau_s} \quad (35)$$

with $\tau_s = .3$ of the passive membrane time constant for all synapses. The connection probability from neurons in population A to population B was p^{AB} , with $p^{EE} = .2$ and $p^{EI} = p^{IE} = p^{II} = .4$. Synaptic weights for connections between excitatory neurons were $\mathbf{W}^{EE} = .0075$ and $\mathbf{W}^{IE} = .0037$, $\mathbf{W}^{EI} = -.0375$, $\mathbf{W}^{II} = -.0375$. These parameters, and the bias voltages μ , were chosen so that the mean field theory derived above was valid for the spiking network's firing rates.

1063 The excitatory neurons were divided into four clusters, each excitatory neuron receiving
1064 half of its inputs from neurons in the same cluster and half from others. Projections to and
1065 from inhibitory neurons were unclustered.

1066 External input from outside the network was contained in I_i^{ext} . We modeled this as a
1067 partially correlated Gaussian white noise process: $I_i^{\text{ext}}(t) = \sigma_i (\sqrt{1-c}\xi_i(t) + \sqrt{c}\xi_c(t))$. $\xi_i(t)$
1068 was Gaussian white noise private to neuron i and $\xi_c(t)$ was shared between all neurons. $c = .05$
1069 denoted the fraction of common input and the noise intensity for excitatory neurons was
1070 $\sigma_E = .3$ and for inhibitory neurons $\sigma_I = .35$.

1071 The firing rate of neuron i in a trial of length L is given by its spike count in that trial n_i^L ,
1072 $r_i = \langle n_i^L \rangle / L$ where $\langle \cdot \rangle$ denotes averaging over trials. The spike train covariance between neurons
1073 i and j describes the above-change likelihood that action potentials occur in each spike train
1074 separated by a time lag s :

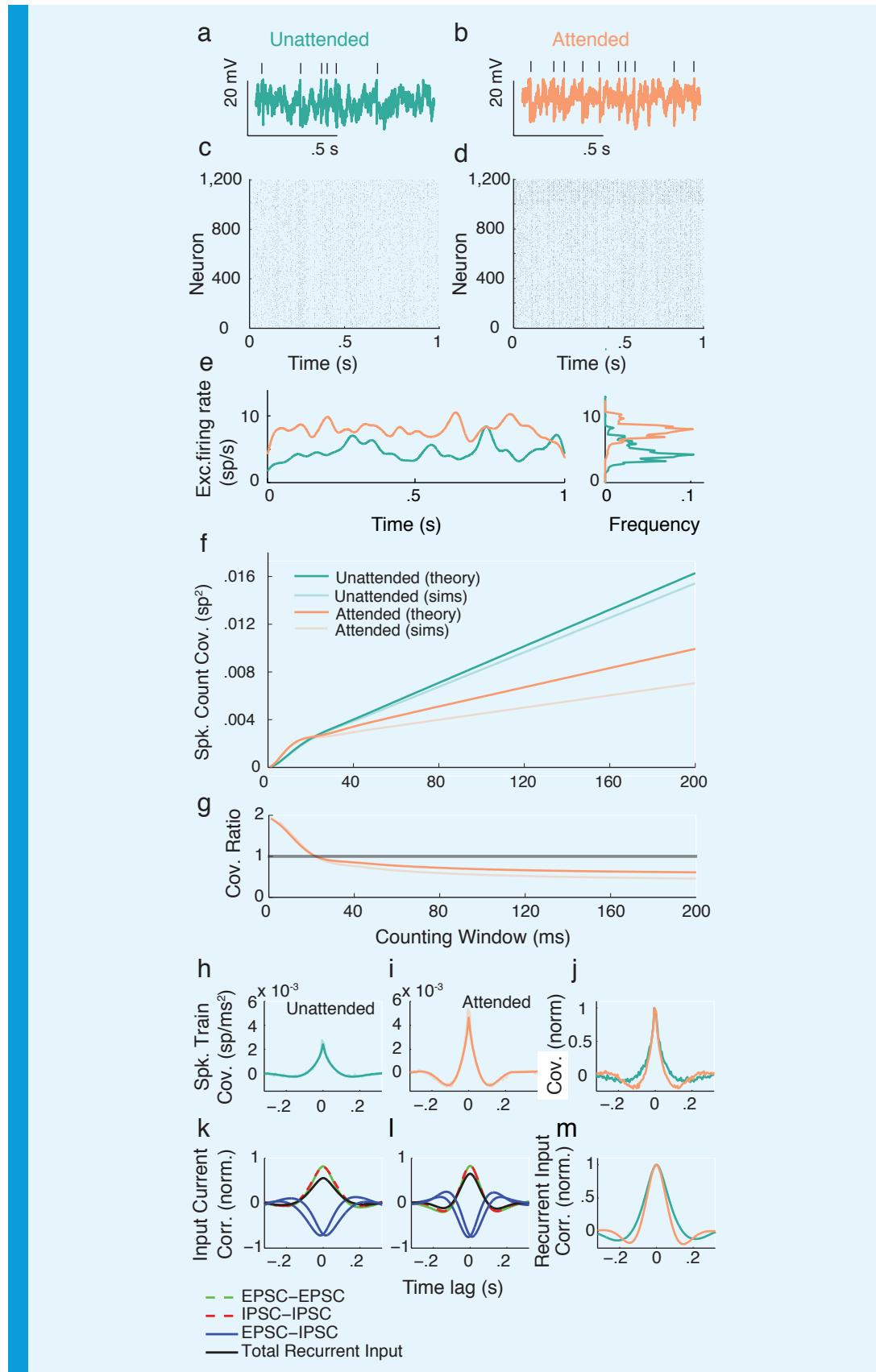
$$\mathbf{1075} \quad q_{ij}(s) = \frac{1}{L} \int_0^L \langle y_i(t) y_j(t-s) \rangle dt - r_i r_j. \quad (36)$$

1076 For simulations, we measure the population-averaged spike train cross-covariance function
1077 $Q(s) = (N_E(N_E - 1))^{-1} \sum_{i,j=1, i \neq j}^{N_E} q_{ij}(s)$ by average a randomly chosen subsample of 100 spike
1078 train cross-covariances from pairs of neurons in the same cluster.

1082 In order to calculate the covariance of neuron i and j 's spike counts in windows of length T ,
1083 n_i^T and n_j^T , we use the relation

$$\mathbf{1085} \quad \text{Cov} \left(n_i^T, n_j^T \right) \equiv \langle n_i^T n_j^T \rangle - \langle n_i^T \rangle \langle n_j^T \rangle \\ \mathbf{1086} \quad = \int_{-T}^T q_{ij}(s) (T - |s|) \\ \mathbf{1087} \quad$$

1088 The cross-correlation of input currents was averaged over the same random subsample of the
1089 network as the spike train covariances. Current cross-correlations were normalized so that each
1090 current's autocorrelation at zero lag was 1.



1092

1093

Appendix 1 Figure 6. Caption is on next page.

Appendix 1 Figure 6. Spiking model and simulations. **a,b** Example voltage trace from an excitatory model neuron in the unattended (a) and attended (b) states. Top tick marks denote spike times. **c**, Raster plot of neurons in the unattended state. Neurons 1 to 1000 are excitatory, and 1001 to 1200 are inhibitory. **d**, Raster plot of neurons in the attended state. **e**, Excitatory population-averaged firing rates for the unattended (turquoise) and attended (orange) states. Right: frequency distributions of population-averaged firing rates. **f**, Mean pairwise spike count covariance for different counting windows. Other than an increase in synchrony on very small timescales due to gamma oscillations, the spike count covariance decreases with attention regardless of counting window. **g**, Ratio R_{Cov} of attended and unattended spike count covariance, as a function of counting window. **h, i** Derived (solid turquoise) and simulated (muted turquoise) spike train cross-covariance functions of excitatory neurons in the unattended (h) and attended (i) states, averaged over pairs. **j**, Spike train cross-covariance functions of excitatory neurons in the unattended and attended states, normalized to peak at 1. **k, l** Normalized input current cross-correlation functions of excitatory inputs to pairs of neurons (dashed green), inhibitory inputs to pairs of neurons (dashed red), excitatory and inhibitory inputs to pairs of neurons (blue), and summed excitatory and inhibitory recurrent inputs to pairs of neurons (black), in the unattended (k) and attended (l) states. **m**, Attended (orange) vs unattended (turquoise) recurrent input cross-correlation functions. The excitatory cross-correlation function is narrower, just as for the output cross-covariance function, so the effects are happening on the level of inputs.

1115 Spiking Network Analysis

The LIF model simulates voltages and produces spike trains, from which we can compute firing rates and covariances. Appendix Fig. 6a,b show example voltage traces of individual excitatory neurons, with the spikes they produce shown above. Note that in the attended state, more spikes are produced, corresponding to a higher firing rate. Appendix Fig. 6c,d show rasters for all the neurons in the unattended (c) and attended (d) states. Higher firing rates can be observed, especially for the inhibitory neurons. Averaging the spike trains over the excitatory population gives us the PSTH of the excitatory neurons. Appendix Fig. 6e, left shows the unattended (turquoise) and attended (orange) PSTH smoothed with a sliding Gaussian window with width (std dev) 10 ms. The histograms on the right demonstrate the decrease in population variance with attention.

The spiking model provides the opportunity to directly compute the pairwise spiking covariance, in addition to the population variance. Appendix Fig. 6f shows the pairwise spike count covariance computed over counting windows from 0 to 200 ms. For small counting windows, corresponding to high-frequency correlations, neurons in the attended state have slightly higher spike count covariance. This is consistent with the slightly higher peak in the attended autocovariance function from the mean-field theory (Figure 4e, main text), as well as experimental results (*Fries et al., 2001*). For counting windows greater than 30 ms, the spike count covariance notably decreases with attention. The experiments we are modeling (*Cohen and Maunsell, 2009*) measure spike count correlations over 200 ms counting windows, corresponding to the right-most points in Appendix 6f. The proportional changes in the spike count covariance are expressed in the covariance ratio $R_{\text{Cov}} = \text{Cov}^A(n_1, n_2)/\text{Cov}^U(n_1, n_2)$, shown in Appendix Fig. 6g. Values of R_{Cov} greater than one indicate increased spike count covariance with attention, and values of R_{Cov} less than one indicate decreased spike count covariance with attention. The crossing of the $R_{\text{Cov}} = 1$ line is apparent at counting windows of approximately 30 ms. The theoretical values were computed using linear response theory (*Trousdale et al., 2012*).

To dissect the spike count covariance by different time lags, we consider the spike train covariance function (Eq. 36), which is the pairwise-neuron analogue of the autocovariance function of the population-averaged activity (Fig. ??e, main text). Appendix Fig. 6h,i show the spike train covariance functions of excitatory neurons in the unattended and attended states. To compare the two, Appendix Fig. 6j shows them normalized so that their maximum values are 1. In accordance with our mean-field results, the attended spike train covariance decays faster than the unattended spike train covariance, indicating increased stability in the

attended state.

The spiking model also provides the opportunity to investigate the inputs to individual neurons, something that is difficult to do experimentally, and does not apply to mean-field models. Appendix Fig. 6k,l shows the correlation functions of different types of inputs to a pair of excitatory neurons, averaged over pairs of excitatory neurons, in the unattended (k) and attended (l) states. Computing the correlation functions of the total recurrent input (black curves) reveals that correlations between excitatory inputs (EPSC-EPSC, dashed green), and correlations between inhibitory inputs (IPSC-IPSC, dashed red), are canceled by anti-correlations between excitatory and inhibitory inputs (EPSC-IPSC, blue). This is consistent with the idea of correlation cancellation by inhibitory tracking of excitatory activity (*Renart et al., 2010; Tetzlaff et al., 2012; Ly et al., 2012*). Attention, by shifting the system into a more stable state, allows this cancellation to occur more efficiently, thereby reducing the pairwise covariance. Appendix Fig. 6m shows the input current correlation functions of the total recurrent inputs to pairs of excitatory neurons, normalized to peak at 1. We conclude that the correlation cancellation brought about by recurrent inhibitory feedback suppresses correlations of the total recurrent input, which in turn decreases the output correlations.