

Circuit models of low dimensional shared variability in cortical networks

Chengcheng Huang^{1,2}, Douglas A. Ruff^{2,3}, Ryan Pyle⁴, Robert Rosenbaum^{4,5},
Marlene R. Cohen^{2,3} and Brent Doiron^{1,2*}

¹Department of Mathematics, University of Pittsburgh, Pittsburgh, PA, USA

²Center for the Neural Basis of Cognition, Pittsburgh, PA, USA

³Department of Neuroscience, University of Pittsburgh, Pittsburgh, PA, USA

⁴Department of Applied and Computational Mathematics and Statistics,
University of Notre Dame, Notre Dame, IN, USA

⁵Interdisciplinary Center for Network Science and Applications,
University of Notre Dame, Notre Dame, IN, USA

*To whom correspondence should be addressed; E-mail: bdoiron@pitt.edu.

Abstract

Trial-to-trial variability is a reflection of the circuitry and cellular physiology that make up a neuronal network. A pervasive yet puzzling feature of cortical circuits is that despite their complex wiring, population-wide shared spiking variability is low dimensional with all neurons fluctuating en masse. Previous model cortical networks are at loss to explain this global variability, and rather assume it is from external sources. We show that if the spatial and temporal scales of inhibitory coupling match known physiology, model spiking neurons internally generate low dimensional shared variability that captures the properties of *in vivo* population recordings along the visual pathway. Shifting spatial attention into the receptive field of visual neurons has been shown to reduce low dimensional shared variability within a brain area, yet increase the variability shared between areas. A top-down

14 modulation of inhibitory neurons in our network provides a parsimonious mechanism for this atten-
15 tional modulation, providing support for our theory of cortical variability. Our work provides a crit-
16 ical and previously missing mechanistic link between observed cortical circuit structure and realistic
17 population-wide shared neuronal variability and its modulation.

18 **Introduction**

19 The trial-to-trial variability of neuronal responses gives a critical window into how the circuit structure
20 connecting neurons drives brain activity¹. This idea combined with the widespread use of population
21 recordings has prompted deep interest in how variability is distributed over a population^{2,3}. There
22 has been a proliferation of data sets where the shared variability over a population is low dimen-
23 sional⁴⁻⁹, meaning that neuronal activity waxes and wanes as a group. In accord, one dimensional
24 measures such as local field potentials^{10,11} and summed population firing rates can predict a major-
25 ity of pairwise correlations^{9,12}. Further, the synthesis of diverse population datasets paints a picture
26 where low dimensional shared variability is a signature of cognitive state, such as overall arousal, task
27 engagement and attention^{2,3,13}, as well as predictive of behavioral performance¹⁴. Such low dimen-
28 sional dynamics portend a theory for the genesis and modulation of shared population variability in
29 recurrent cortical networks.

30 Theories of cortical variability can be broadly separated into two categories: ones where vari-
31 ability is internally generated through recurrent network interactions (Fig. 1a, left) and ones where
32 variability originates external to the network (Fig. 1a, middle). Networks of spiking neuron mod-
33 els where strong excitation is balanced by opposing recurrent inhibition produce high single neuron
34 variability through internal mechanisms¹⁵⁻¹⁷. However, these networks famously enforce an asyn-
35 chronous state, and as such fail to explain population-wide shared variability¹⁸. This lack of success
36 is contrasted with the ease of producing arbitrary correlation structure from external sources. Indeed,
37 many past cortical models assume a global fluctuation from an external source^{3,7,19-21}, and accurately

38 capture the structure of population data. However, such phenomenological models are circular, with
39 an assumption of variability from an unobserved source explaining the variability in a recorded pop-
40 ulation. Thus, while neuronal variability has a rich history of study, there remains an impoverished
41 mechanistic understanding of the low dimensional structure of population-wide variability²².

42 Determining whether output variability is internally generated through network interactions or ex-
43 ternally imposed upon a network is a difficult problem, where single area population recordings may
44 preclude any definitive solution (Fig. 1a, left vs middle). In this study we consider attention-mediated
45 shifts in population variability obtained from simultaneous recordings of neuron pairs both within
46 and between visual areas^{23,24}. Attention reduces within area correlations (area V4) while simultane-
47 ously increasing between area correlations (areas V1 and MT), thereby providing a novel constraint
48 for how shared variability is distributed within and between neuronal populations (Fig. 1a, right). We
49 present analysis showing that such a differential correlation modulation is difficult constraint to sat-
50 isfy with a model where fluctuations are strictly external to the network. We thus focus our modeling
51 on networks where population-wide shared variability can be internally generated.

52 The asynchronous solution of classical balanced networks necessitates that inhibition dynami-
53 cally tracks and cancels any correlations steaming from recurrent excitation¹⁸. This requirement has
54 forced theorists to assume that the timecourse of inhibitory synapses is faster than that of excitatory
55 synapses^{16,18,25-27}, at odds with recorded synaptic physiology²⁸. Recently, we have extended the the-
56 ory of balanced networks to include a spatial component to network architecture^{25,26,29} and found
57 network solutions where firing rate balance and asynchronous dynamics are decoupled from one an-
58 other²⁶. In this study, we consider multi-area models of spatially distributed balanced networks and
59 show that when inhibition has slower kinetics than excitation these networks, matching physiology,
60 they internally produce low dimensional population-wide variability. Unlike networks that lack spa-
61 tial structure, these networks produce spiking activity that robustly captures the rich diversity of firing
62 rate and correlated structure of real population recordings. Further, attention-mediated top-down

63 modulation of inhibitory neurons in our model provides a parsimonious mechanism that controls
64 population-wide variability in agreement with the within and between area experimental results.

65 There is a long standing research program aimed at providing a circuit-based understanding for
66 cortical variability^{1,15–17,26}. Our work is a critical advance through providing a mechanistic theory
67 for the genesis, propagation, and modulation of realistic low dimensional population-wide shared
68 variability based on established circuit structure and synaptic physiology.

69 **Results**

70 **Externally imposed or internally generated shared variability?**

71 Directed attention reduces the mean spike count correlation coefficient between neuron pairs in visual
72 area V4 during an orientation detection task (Fig. 1b;²⁴). In V4⁵, as with other cortices^{4,6,8,9,12}, shared
73 variability across a population is low dimensional, where coordinated fluctuations are driven by a
74 common latent variable. Further, attention reduces pairwise correlation through attenuation of this
75 global latent variable^{5,30}. Thus motivated, we represent the aggregate population response with a
76 scalar random variable $R = X + \beta H$, where X is a noisy stimulus input and H is a hidden source
77 of fluctuations (with strength β ; Fig. 1c, top). In this simple model the trial-to-trial fluctuations are
78 inherited from both X and H , but we model attention as only reducing the variance of H ($\text{Var}(H)$).
79 There is a large range of parameter values for our one-dimensional hidden variable model to readily
80 explain the reduction in $\text{Var}(R)$ reported in the V4 data (Fig. 1c, bottom, blue curve; see Supplemental
81 Information). Certain parameter choices are unreasonable (pink region in Fig. 1c, bottom), such as β
82 being overly large so that R is no longer driven by X , or $\text{Var}(H) \rightarrow 0$ in the attended state, requiring
83 the area that produces H to be silent. Fortunately, there are moderate β and $\text{Var}(H)$ choices that
84 capture the data (section of the blue curve that is not in the pink region in Fig. 1c, bottom). In total, a
85 latent variable model where the variability is external to the population can account for the attentional
86 modulation reported in our V4 data, as has been previously remarked^{5,7}.

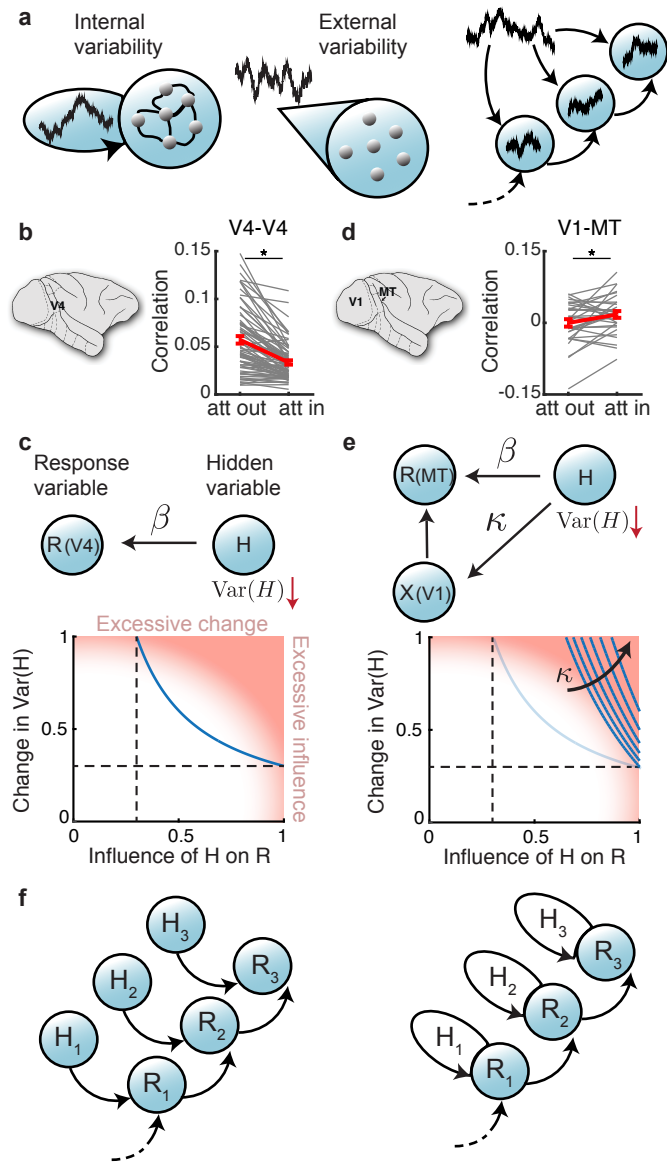


Figure 1: Caption is on next page.

Figure 1: **Models of shared variability.** **a**, Variability may either be internally generated within a population (left) or externally imposed upon a population (middle). New model constraints emerge by accounting how variability is distributed and modulated across several populations (right). **b**, Mean spike count correlation r_{SC} per session obtained from multi-electrode array recording from V4 was smaller when attention was directed into the receptive fields of recorded neurons (n=74 sessions, two-sided Wilcoxon rank-sum test between attentional states $P = 3.3 \times 10^{-6}$, reproduced from²⁴). Grey lines are individual session comparisons and the red line is the mean comparison across all sessions (error bars represent the SEM). **c**, Top: hidden variable model where the response variability R (modeling V4) comes from a hidden variable H with influence β . Bottom: the attention-mediated reduction in r_{SC} gives a constraint that is a trade-off between the reduction in $\text{Var}(H)$ and β (blue curve). **d**, Same as **b** for the mean spike count correlation r_{SC} between V1 units and MT units per session (n=32 sessions, paired-sample t -test $P = 0.0222$; data reproduced from²³). **e**, Top: hidden variable model for connected areas X (modeling V1) and R (modeling MT); H projects to X with strength κ . Bottom: the attention mediated changes in r_{SC} give further constraints on H with the increase in κ indicated. Light blue curve is the same as that in **c** for comparison. **f**, Schematic for external (left) and internal (right) models of shared variability (H) along a processing hierarchy (R).

87 Recent multi-electrode recordings from two visual areas, MT and V1, during an attention task²³
88 impose strong constraints on the simple hidden variable model. In addition to a reduction of mean
89 spike count correlations between neuron pairs within an area (pairwise attention-related MT correla-
90 tion decrease, 0.019, Wilcoxon rank sum test, $p = 0.017$; pairwise attention-related V1 correlation
91 decrease, 0.008, Wilcoxon rank sum test, $p = 4.9 \times 10^{-6}$;²³), there is an attention-mediated *increase*
92 of spike count correlations across areas V1 and MT (Fig. 1d). Returning to the population model
93 with R modeling MT, we augment the model with V1 being the input $X = X_0 + \kappa H$ (Fig. 1e, top;
94 see Supplementary Information). Here κ denotes how much the hidden variable is directly shared
95 between areas, and X_0 is the variability in X that is independent of H . The constraint curves (Fig.
96 1e, bottom blue) where $\text{Var}(R)$ and $\text{Cov}(R, X)$ match the MT-MT and V1-MT data sets require our
97 model to assume both a large influence of H on R and a large attentional modulation of $\text{Var}(H)$ (pink
98 region in Fig. 1e, bottom). This tightening of model assumptions reflects the compromise between
99 an attention-mediated increase in variability transfer from $X \rightarrow R$ so that $\text{Cov}(R, X)$ increases and
100 a simultaneous decrease in $\text{Var}(H)$ so that $\text{Var}(R)$ decreases. This compromise can be mitigated by

101 setting κ to be small, meaning that a large component of the fluctuations in R is private from those in
102 X (Fig. 1e, bottom).

103 While the source of private variability H to area R may still be external to the area, if we extrap-
104 olate our model to a cortical hierarchy then each area requires an external variability ‘generator’ that
105 projects privately to that area (Fig. 1f, left). This would require a tremendous amount of neuronal
106 hardware. A more parsimonious hypothesis is that private variability is internally generated within
107 each area (Fig. 1f, right). Below we investigate the circuit mechanics required for low dimensional
108 population-wide shared variability to be an emergent property within a cortical network.

109 **Population-wide correlations with slow inhibition in spatially ordered networks**

110 Networks of spiking neuron models where strong excitation is balanced by opposing recurrent inhi-
111 bition internally produce high single neuron variability (Fig. 2ai) with a broad distribution of firing
112 rates (Fig. 2b, top purple curve)^{16–18}. However, these networks enforce an asynchronous solution
113 (Fig. 2c, top purple), and as such fail to explain population-wide shared variability^{8,18}. Typically,
114 balanced networks have disordered connectivity, namely where connection probability is uniform be-
115 tween all neuron pairs. This approximation ignores the abundant evidence that cortical connectivity is
116 spatially ordered with a connection probability falling off with the distance between neuron pairs^{31–33}.
117 Recently we have extended the theory of balanced networks to include such spatially dependent con-
118 nectivity^{25,26}. Briefly, we model a two dimensional array of integrate-and-fire neurons that receive
119 both feedforward projections from a layer of external Poisson processes and recurrent projections
120 within the network (see Methods); connection probability of all projections decays like a Gaussian
121 with distance. If the spatial scale of feedforward inputs is narrower than the scale of recurrent pro-
122 jections, the asynchronous state no longer exists²⁶, giving way to a solution with spatially structured
123 correlations (Fig. 2aii, Supplemental movie S1, Fig. S1b). Nevertheless, the mean correlation across
124 all neuron pairs vanishes for large network size (Fig. 2c, bottom purple curve), in stark disagreement

125 with a majority of experimental studies^{2,3} (Fig. 1b,d).

126 Many previous balanced network models assume that the kinetics of inhibitory conductances are
127 *faster* than those of excitatory conductances^{16–18,26,34}. However, this assumption is at odds with phys-
128 iology where excitatory α -amino-3-hydroxy-5-methyl-4-isoxazolepropionic acid receptors (AMPA)
129 have faster kinetics than those of the inhibitory γ -Aminobutyric acid receptors (GABA_A)²⁸. When
130 the timescales of excitation and inhibition match experimental values in networks with disordered
131 connectivity the activity becomes pathologic, with homogeneous firing rates (Fig. 2b, top green) and
132 excessive synchrony (Fig. 2a_{iii}, and c, top green). This consequence is likely the ad-hoc justification
133 for the faster inhibitory kinetics in disordered model networks.

134 When a spatially ordered model has synaptic kinetics that match physiology, population-wide tur-
135 bulent dynamics emerges (Fig. 2a_{iv}, Supplemental movie S2), accompanying a small, but nonzero,
136 mean pairwise spike count correlation across the population ($r_{SC} = 0.04$). Further, firing rates are
137 broad (Fig. 2b, bottom green curve) and pairwise correlations are reasonable in magnitude (Fig. 2c,
138 bottom green). Indeed, as the timescale of inhibition grows, disordered networks show a rapid change
139 in mean pairwise correlation while two dimensional spatially ordered networks show a much more
140 gradual rise in correlation (Fig. 2d). We remark that networks constrained to one spatial dimen-
141 sion also produce excessive synchrony (Fig. S2), meaning that two (or more) spatial dimensions are
142 required for robustly low but nonzero correlations. In sum, when realistic spatial synaptic connec-
143 tivity is paired with realistic temporal synaptic kinetics in balanced networks, internally generated
144 population dynamics produces spiking dynamics whose marginal and pairwise variability conform to
145 experimental results.

146 **Attentional modulation of low dimensional population-wide variability**

147 We model the V1 and MT network by extending the spatially ordered balanced networks with slow
148 inhibition to include three layers: a bottom layer of independent Poisson processes modeling thala-

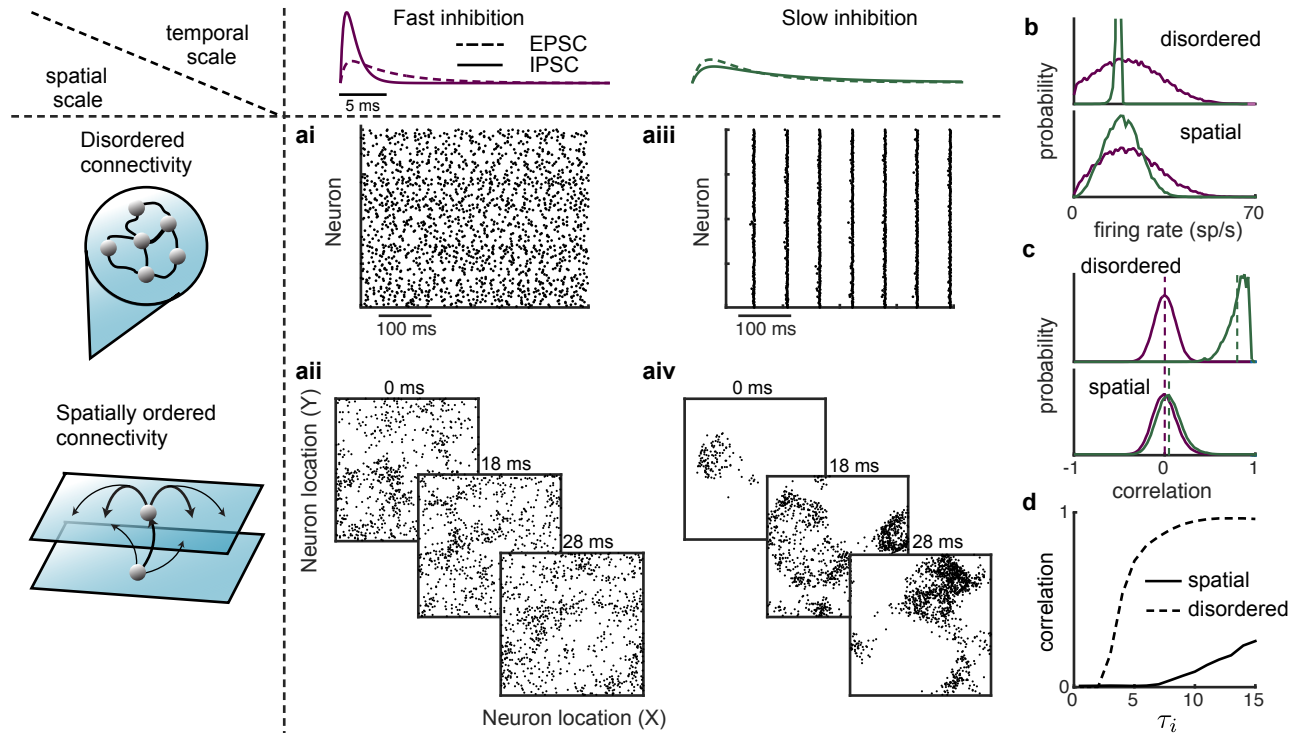


Figure 2: The spatial and temporal scales of synaptic coupling determine internally generated variability. **a**, Networks of excitatory and inhibitory neuron models were simulated with either disordered connectivity (ai,aiii) or spatially ordered connectivity (aii,aiv), and with either fast inhibition ($\tau_i = 1$ ms; ai,aii). or slow inhibition ($\tau_i = 8$ ms; aiii,aiv). In all models the timescale of excitation was $\tau_e = 5$ ms. In the disordered networks spike train rasters assume no particular neuron ordering. In the spatially order networks three consecutive spike raster snapshots are shown with a dot indicating that the neuron at spatial position (x, y) fired within one millisecond of the time stamp. **b**, Distributions of firing rates of excitatory neurons in the disordered (top) and spatially ordered (bottom) models, with faster inhibitory kinetics (purple) compared to slower inhibitory kinetics (green). **c**, Same as **b** for the distributions of pairwise correlations among the excitatory population. **d**, Mean correlation among the excitatory population as a function of the inhibitory decay time constant (τ_i).

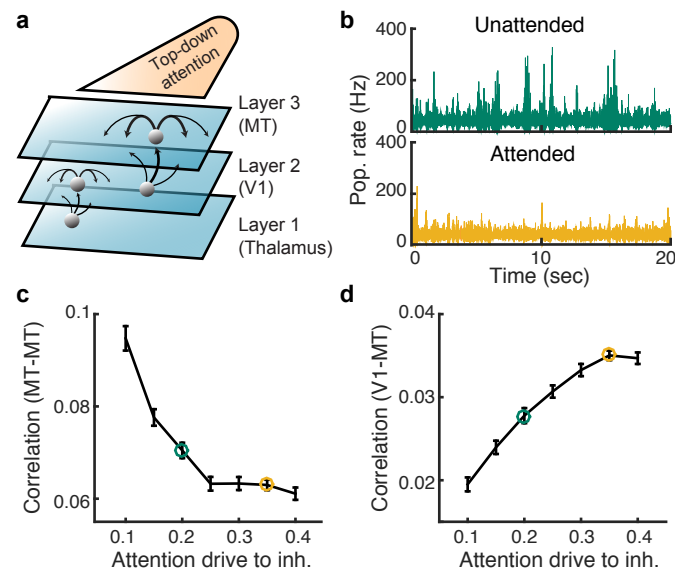


Figure 3: **Top-down depolarization of MT inhibitory neurons capture the differential attentional modulation of shared variability within and across V1 and MT.** **a**, Thalamus, V1, and MT are modeled in a three layer hierarchy of spatially ordered balanced networks. Top-down attentional modulation is modeled as a depolarization to MT inhibitory neurons (μ_I). In both V1 and MT the recurrent projections are broader than feedforward projections and recurrent inhibition is slower than excitation. **b**, Population averaged firing rate fluctuations from MT in the unattended state ($\mu_I = 0.2$, green) and the attended state ($\mu_I = 0.35$, orange). **c**, Mean spike count correlation (r_{SC}) of excitatory neuron pairs in MT decreases with attentional modulation. **d**, Mean r_{SC} between the excitatory neurons in MT and the excitatory neurons in V1 increases with attention. Error bars are SEM.

149 mus, and middle and top layers of integrate-and-fire neurons modeling V1 and MT, respectively (Fig.
150 3a and see Methods). We follow our past work with simplified firing rate networks⁷ and model a
151 top-down attentional signal as an overall static depolarization to inhibitory neurons in the MT layer
152 (Fig. 3a). This mimics cholinergic pathways that primarily affect interneurons^{35,36} and are thought
153 to be engaged during attention^{7,13}. The increased recruitment of inhibition during attention reduces
154 the population-wide fluctuations in the MT layer (Fig. 3b) and decreases pairwise spike count corre-
155 lations of MT-MT neuron pairs (Fig. 3c), while simultaneously increasing the correlation of V1-MT
156 neuron pairs (Fig. 3d). Thus, this simple implementation of attentional modulation⁷ nonetheless cap-
157 tures the main aspects of the V1-MT dataset (Fig. 1d;²³). Further, neuron pairs with larger firing rate
158 increases also show larger correlation reductions (Fig. S3), in agreement with population recordings
159 during both spatial and feature attention³⁷.

160 Before we expose the core mechanisms through which attention modulates correlated activity we
161 first give a broader analysis of shared variability in both our data and model. To this end we use
162 dimensionality reduction tools^{8,38} to study the population-wide structure of trial-to-trial variability,
163 rather than focusing only on individual pairwise correlation coefficients. We partition the covariance
164 matrix into the shared variability among the population and the private noise to each neuron; the
165 eigenvalues of the shared covariance matrix represent the variance along each dimension (or latent
166 variable), while the corresponding eigenvectors represent the projection weights of the latent variables
167 onto each neuron (see Methods). Applying these techniques to the multi-electrode V4 data²⁴ shows
168 a single dominant eigenmode (Fig. 4a, top, single session result see Fig. S4). This mode influences
169 most of the neurons in the population in the same way (Fig. 4a middle, weights are dominant positive),
170 and after subtracting the first mode the mean residual covariances are very small (Fig. 4a bottom).
171 Finally, attention affects population variability primarily by quenching this dominant mode (Fig. 4a
172 top, orange vs green) and the attentional modulation in the dominant mode is highly correlated with
173 the modulation in mean covariance (Fig. S4c). The low dimensional structure of shared variability in

174 our data is consistent with similar analysis in other cortices^{4,6,8}, as well as alternative analysis of the
175 same V4 data using generalized point process models⁵.

176 The dimensionality of shared variability offers a strong test for our cortical model. We analyzed
177 the spike count covariance matrix constructed from a subsampling of the spike trains in the third layer
178 of our network model ($n = 50$ neurons). The network with slow inhibition produced shared variability
179 with a clear dominant eigenmode that mimicked many of the core features observed in the V4 data
180 (Fig. 4b). Further, the top-down attentional modulation of inhibition also suppressed this dominant
181 mode (Fig. 4b top, orange vs green). The agreement between model and data broke down when
182 inhibitory kinetics were faster than those of excitation, as was the case in our past studies^{25,26,29}. Here,
183 shared variability did not have a dominant mode (Fig. 4c, top), the raw mean correlation coefficient
184 was near zero (Fig. 4c, bottom), and attentional modulation had a negligible effect on population
185 variability (Fig. 4c, orange vs green). Experimental measurements of local cortical circuitry show that
186 excitation and inhibition project on similar spatial scales^{31,33}. When the model inhibitory projections
187 in the third layer were spatially broader than those of excitation, thus at odds with experiment, then the
188 model again disagreed with our V4 data (Fig. 4d). In sum, the low dimensional structure of shared
189 variability requires inhibition that is neither faster nor anatomically broader than excitation – both
190 features of real cortical circuits^{28,31,39}. Further, a simple recruitment of inhibition through top-down
191 drive can restore stability and quench low dimensional population variability.

192 This success of our model is quite distinct from that of past studies where low dimensional corre-
193 lated variability was imposed from outside sources^{3,7,19–21}. Rather, the shared variability in our model
194 is internally generated from recurrent network interactions. We next explore how the inherent nonlin-
195 ear dynamics that produce this variability allow our model to satisfy the constraints imposed by the
196 differential correlation modulation of the within area and between area pairs (Fig. 1e).

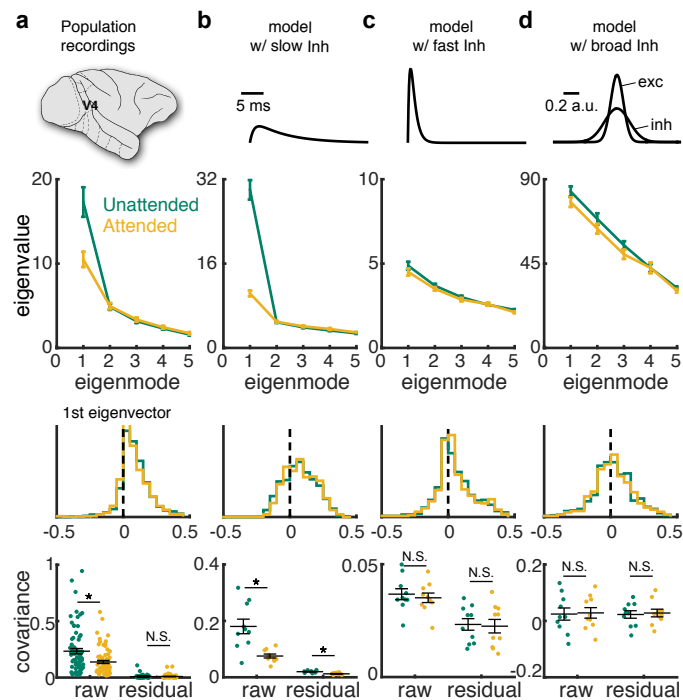


Figure 4: Internally generated shared variability from the model network is low-dimensional.

a, Top: The first five largest eigenvalues of the shared component of the spike count covariance matrix from the V4 data²⁴. Green: unattended; orange: attended; data from $n=72$ sessions with 43 ± 15 neurons. Error bars are SEM. Middle: the vector elements for the first (dominant) eigenmode. Bottom: the mean covariance from each session in attended and unattended states before (raw) and after (residual) subtracting the first eigenmode (mean \pm s.e.m. in black). **b-d**, Same as **a** but for the three layer model with slow inhibition (**b**), model with fast inhibition (**c**) and model with slow and broad inhibition (**d**); $n=10$ samples of 50 neurons each. Two-sided Wilcoxon rank-sum test (attended vs unattended): mean covariance, **a**, $P = 0.0013$, **b**, $P = 1.78 \times 10^{-22}$, **c**, $P = 0.7798$ and **d**, $P = 0.5850$; residual, **a**, $P = 0.7477$, **b**, $P = 5.40 \times 10^{-4}$, **c**, $P = 0.8796$ and **d**, $P = 0.5326$.

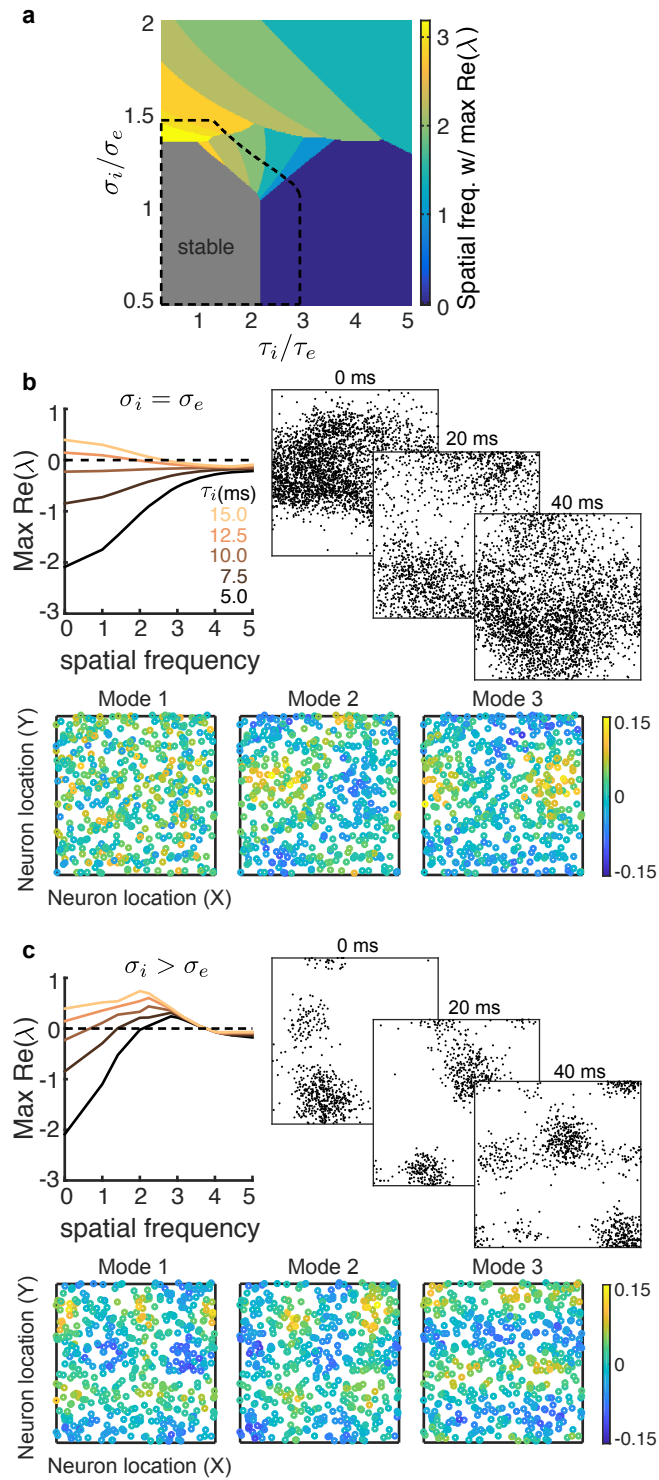


Figure 5: Caption is on next page.

Figure 5: **Stability analysis of a two-dimensional firing rate model.** **a**, Bifurcation diagram of a firing rate model as a function of the inhibitory decay time scale τ_i and inhibitory projection width σ_i . The excitatory projection width and time constant are fixed at $\sigma_e = 0.1$ and $\tau_e = 5$ ms, respectively. Color represents the spatial frequency with the largest real part of eigenvalue and the gray region is stable. Top-down modulation of inhibitory neurons modeling attention expands the stable region (black dashed). **b**, Top left: the real part of eigenvalues as a function of spatial frequency for increasing τ_i when $\sigma_i = \sigma_e$. Top right: three consecutive spike raster snapshots of a spiking neuron network with $\sigma_i = \sigma_e$ and slow inhibition (same network as in Fig. 4b in the unattended state). Bottom: spatial structure of projection weights from the first three eigenmode from Factor analysis of the spiking neuron network as in top right ($n=500$ neurons). **c**, Same as **b** for σ_i larger than σ_e . Top right and Bottom: same network as in Fig. 4d in the unattended state.

197 **Relating low dimensional variability to spatio-temporal pattern formation**

198 Networks of spiking neuron models produce rich activity that can be directly compared to population
199 recordings. However, when these networks are outside the asynchronous regime they are not easily
200 amenable to a deeper mechanistic analysis. An often used simplification to spiking dynamics are fir-
201 ing rate models where network interactions are mediated only through dynamic firing rates^{1,40}. While
202 these models lack a principled connection to spiking network models, they do produce qualitatively
203 similar dynamics in recurrent networks and their simplicity makes them amenable to analysis tech-
204 niques from dynamical systems theory. To gain intuition about how recurrent circuitry shapes low
205 dimensional shared variability we considered a firing rate model that incorporated both the spatial
206 architecture and synaptic dynamics that were central to our spiking model (see Methods).

207 Solutions where firing rates are constant over time are interpreted as asynchrony within the net-
208 work, since only dynamical co-fluctuations in firing rates would mimic correlated spiking. We fo-
209 cused on how the stability of the asynchronous firing rate solution depended upon the temporal (τ_i)
210 and spatial (σ_i) scales of inhibition. A firing rate solution is stable if the linearized dynamics are such
211 that every eigenmode has eigenvalues with strictly negative real part. Since our network is spatially
212 ordered the eignemodes are also organized in space, each with their own distinct spatial frequency. If
213 the solution loses stability at a particular eigenmode, then the spatio-temporal dynamics of the result-

214 ing network firing rates will inherit the spatial frequency of that eigenmode – this process is termed
215 spatio-temporal pattern formation⁴¹.

216 If τ_i and σ_i are near those of recurrent excitation, then a stable firing rate solution exists (Fig.
217 5a, grey region; 5b, top left, black curve with $\tau_i = 5\text{ms}$). Our past work explored activity within
218 this regime²⁶. When τ_i increases and excitation and inhibition project with the same spatial scale
219 ($\sigma_i = \sigma_e$), firing rate stability is first lost at an eigenmode with zero spatial frequency (Fig. 5b, top
220 left). This creates population dynamics with a broad spatial pattern, allowing variability to be shared
221 over the entire network. Simulations of the three layered spiking network model in this regime shows
222 turbulent dynamics that extend across the entire network (Fig. 5b, top right; Supplementary movie
223 S3). The projection weights of the first eigenmode from factor analysis of the network of spiking
224 neuron models (Fig. 4b) show a uniform distribution in space (Fig. 5b, bottom), consistent with
225 shared fluctuations of low spatial frequency. In contrast to this case, when τ_i increases yet inhibition
226 projects lateral to excitation ($\sigma_i > \sigma_e$), stability is first lost at a nonzero spatial frequency (Fig. 5c,
227 top left). This creates population dynamics with coherence over a band of higher spatial frequencies,
228 producing higher dimensional shared variability, as evident in the spatially patchy turbulent spiking
229 dynamics of the three layered spiking network in this regime (Fig. 5c, top right; Supplementary movie
230 S4). Correspondingly, the projection weights of the first three eigenmodes (Fig. 4d) show patterns
231 of higher spatial frequency (Fig. 5c, bottom). Thus, the spatial and temporal scales of inhibition
232 determine in large part the spatio-temporal patterns of network activity. Further, we now understand
233 from a theoretical viewpoint why slow inhibition that does not project lateral to excitation is needed
234 to account for the spiking data in both the network of spiking neuron models and experiment.

235 Finally, in the firing rate network we can also model attention as a depolarization to the inhibitory
236 neurons, as was done in the network of spiking neuron models. In the firing rate network, attentional
237 modulation expanded the stable region in the bifurcation diagram (Fig. 5a, dashed black line). In
238 other words, attention increased the domain of firing rate stability. Thus, with $\tau_i > \tau_e$ chosen so

239 that in the unattended state the network was unstable at a low spatial frequency yet with attention
240 the network was in the stable regime, our model captures the large attention-mediated quenching of
241 population-wide shared variability reported in the population recordings (Fig. 4a) and network of
242 spiking neuron models (Fig. 4b).

243 **Chaotic population-wide dynamics reflects internally generated variability**

244 The attention-mediated differential modulation of within and between area correlations lead us to
245 propose that shared variability has a sizable internally generated component. Using our heuristic
246 model we argued that attention must quench a significant component of the variability ($\text{Var}(H)$) to
247 account for the population recordings (Fig. 1e, bottom). This is a difficult constraint to satisfy and
248 requires the mechanisms that produce internally generated variability to sensitively depend on top-
249 down modulations. The firing rate model captured this sensitivity through a spatio-temporal pattern
250 forming transition in network activity. However, the firing rate model does not internally produce trial-
251 to-trial variability that can be compared to experiment, and we thus return to analysis of the network of
252 spiking neuron models to probe how trial-to-trial variability is internally generated through recurrent
253 coupling.

254 To isolate the sources of externally and internally generated fluctuations in the third layer of our
255 network we fixed the spike train realizations from the first layer (thalamic) neurons as well as the
256 membrane potential states of the second layer (V1) neurons, and only the initial membrane potentials
257 of the third layer (MT) neurons were randomized across trials (Fig. 6a). This produced deterministic
258 network dynamics when conditioned on activity from the first two layers, and consequently any trial-
259 to-trial variability is due to mechanics internal to the third layer.

260 The spike trains from third layer neurons in both the unattended and attended states have signif-
261 icant trial-to-trial variability despite the frozen layer one and two inputs. This is reflective of a well
262 studied chaotic network dynamic in balanced networks where the spike times from individual neu-

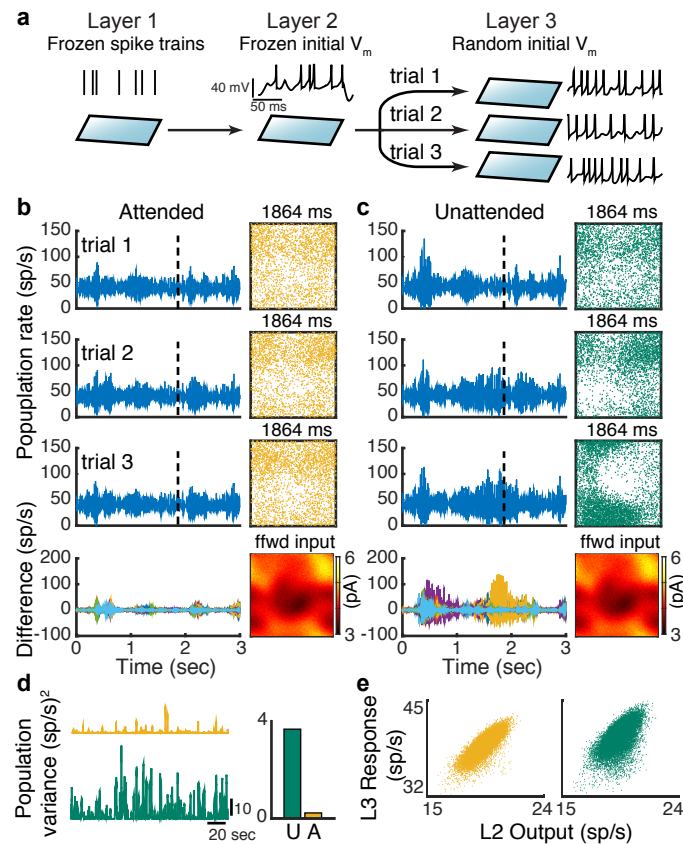


Figure 6: Chaotic population firing rate dynamics is quenched by attention. **a**, Schematic of the numerical experiment. The spike train realizations in layer one and the initial states of the membrane potential of layer two neurons are identical across trials, while in each trial we randomized the initial states of the layer three neuron's membrane potentials. **b**, Three representative trials of the layer three excitatory population rates in the attended state (left row 1-3). Bottom row: difference of the population rates across 20 trials. Right (row 1-3): Snapshots of the neuron activity at time point 1864 ms. Each dot is a spike within 2 ms window from the neuron at that location. Right bottom: the synaptic current each layer three neuron receives from layer two at time 1864 ms. **c**, Same as **b** for the network in the unattended state. **d**, Trial-to-trial variance of layer three population rates as a function of time; right: mean variance across time. **e**, The layer three population rate tracks the layer two population rate better in the attended state. Both outputs and responses are smoothed with a 200 ms window.

263 rons are very sensitive to perturbations that affect the spiking of other neurons^{27,42}. To investigate how
264 this microscopic (single neuron) variability possibly manifests as macroscopic population activity, we
265 considered the trial-to-trial variability of the population-averaged instantaneous firing rate. While the
266 population firing rate is dynamic in the attended state, there is very little variability from trial-to-trial
267 (Fig. 6b, left; Fig. 6d, orange). A consequence of this low population-wide variability is the faithful
268 tracking of the spatiotemporal structure of layer two outputs by layer three responses (Fig. 6b, right;
269 e, orange). This tracking reflects the higher correlation between layer two and three spiking in the
270 attended state (Fig. 3c). In contrast, in the unattended state the asynchronous solution is unstable,
271 resulting in population-wide recruited activity. These periods of spatial coherence across the network
272 are not trial locked and rather contribute to sizable trial-to-trial variability of population activity (Fig.
273 6c, left; d, green). This degrades the tracking of layer two outputs (Fig. 6c, right; e, green) and
274 ultimately lowers the correlation between layer two and three spiking (Fig. 3c). Taken together, while
275 the network model is chaotic in both the attended and unattended states, the chaos is population-wide
276 only in the inhibition deprived unattended state.

277 The nonlinear pattern forming dynamics of the spatially distributed recurrent network impart ex-
278 treme sensitivity to the population-wide internally generated variability. Indeed, in our model the
279 trial-to-trial population rate variability is almost extinguished with attention (Fig. 6d, right). In our
280 heuristic model with hidden variable H this amounts to $\text{Var}(H)$ reducing drastically with attention,
281 which is precisely what is needed to account for the differential modulation of within and between
282 area correlations (Fig. 1e).

283 Discussion

284 There is a longstanding research program aimed at understanding how variability is an emergent prop-
285 erty of recurrent networks^{16,17,25–27,42}. However, models are often restricted to simple networks with
286 disordered connectivity. Consequently, population-wide activity is asynchronous, at odds with many

287 experimental findings^{2,3}. A parallel stream of research focuses on spatiotemporal pattern formation
288 in neuronal populations, with a rich history in both theoretical⁴⁰ and experimental contexts⁴³. Yet a
289 majority of these studies consider only trial-averaged activity, with tacit assumptions about how spik-
290 ing variability emerges (but see⁴⁴ and²⁶). In this study we combined these modelling traditions with
291 the goal of circuit-based understanding of the genesis and modulation of low dimensional internally
292 generated shared cortical variability.

293 **Population-wide variability in balanced networks**

294 Our model extends classical work in balanced cortical networks^{16,18} to include two well accepted ex-
295 perimental observations. First, cortical connectivity has a wiring rule that depends upon the distance
296 between neuron pairs^{31,32}. Theoretical studies that model distance dependent coupling commonly as-
297 sume that inhibition projects more broadly than excitation^{40,44,45} (but see³⁴). However, measurements
298 of local cortical circuitry show that excitation and inhibition project on similar spatial scales^{31,33}, and
299 long-range excitation is known to project more broadly than inhibition⁴⁶. Our work shows that this
300 architecture is required for internally generated population variability to be low dimensional (Fig.
301 4b, d). The second observation is that inhibition has temporal kinetics that are slower than excita-
302 tion²⁸. Past theoretical models of recurrent cortical circuits have assumed that inhibition is not slower
303 than excitation^{16,18,34,47}, including past work from our group²⁶. Consequently, these studies could
304 only capture the residual correlation structure of population recordings once the dominant eigenmode
305 was subtracted^{8,26}; in these cases the residual accounted for less than ten percent of the true shared
306 variability. The asynchronous solution is unstable when inhibition is slower than excitation, and in
307 networks with two spatial dimensions the resulting dynamics are weakly correlated, matching experi-
308 ments (Fig. 2 and 4). In total, by including accepted features of cortical anatomy and physiology, long
309 ignored by theorists, our model network recapitulates low dimensional population-wide variability to
310 a much larger extent than previous models.

311 The above narrative is somewhat revisionist; there are several well known theoretical studies
312 in disordered networks where one dimensional population-wide correlations do emerge, notably in
313 networks where rhythmic¹⁷ or ‘up-down’^{45,47} dynamics are prominent. Networks with dense yet dis-
314 ordered connectivity ensure that all neuron pairs receive some shared inputs from overlapping presy-
315 naptic projections. In such a network if the asynchronous state becomes unstable then this shared
316 wiring will correlate spiking activity across the entire network. In other words, any shared variability
317 will be one dimensional (scalar) by construction. In contrast, the ordered connectivity in our network
318 is such that neuron pairs that are distant from one another have no directly shared presynaptic con-
319 nections. Consequently, when asynchrony is unstable one dimensional population dynamics is not
320 preordained, rather the spatial network can support higher dimensional shared variability depending
321 on the temporal and spatial scales of recurrent coupling (Fig. 4b,c; Fig. 5). From the vantage of
322 this model we discovered the conditions for recurrent architecture and synaptic physiology for low
323 dimensional shared variability

324 **Internal versus external population variability**

325 Our circuit model assumed that the component of population-wide variability that is subject to atten-
326 tional modulation was internally generated within the network. This was motivated by constraints
327 imposed by the differential attentional modulation of within and between pairwise correlations in our
328 population recordings²³ (Fig. 1). While our model is a parsimonious explanation of the data, it does
329 not definitively exclude mechanisms where variability is inherited from outside sources. In fact it is
330 difficult to conceive of descending synaptic and cholinergic projections from higher areas that would
331 not contribute some trial-to-trial variability to a receiving neuronal population.

332 Fluctuations from external sources are an often assumed and straightforward mechanism for
333 population-wide variability^{3,7,19–21,48}. However, if this framework aims to capture a modulation in
334 variability a further choice must be made³. One way to modulate population-wide variability is to

335 simply allow the amplitude of input fluctuations to change. Such an ‘inheritance model’ is often as-
336 sumed for how top-down feedback to either visual areas V1⁴⁸ or MT²¹ determines choice probability
337 in ambiguous decision tasks. When V2 and V3 are inactivated through cooling the single neuron
338 variability in MT is markedly reduced suggestive that a component of variability is feedforward prop-
339 agated⁴⁹. This is in contrast to the only slight reductions in V1 variability when feedback projections
340 from V2 and V3 are inactivated⁴⁹. Thus, there is limited experimental evidence for direct top-down
341 contributions to single neuron variability. Additional multi-area population recordings between con-
342 nected brain regions will be needed to probe how correlated variability flows along bottom-up and
343 top-down pathways.

344 The second way to change population output variability is to keep input fluctuations fixed yet
345 shift the operating point of the network so that the nonlinearities inherent in spiking dynamics change
346 input-output transfer of variability. This mechanism has been suggested for how top-down attentional
347 modulation affects population variability in recurrent excitatory-inhibitory cortical networks^{7,19}. Net-
348 work models with either disordered connectivity or simple one dimensional spatial structure must
349 have a stable asynchronous state, else the internally generated correlations are excessive (Fig. 2aiii,c).
350 Consequently, when such networks are used to model attentional modulation both the attended and
351 unattended states must be in the asynchronous regime^{7,19}. In such cases, population-wide variability
352 must be from outside the network and attention only changes how the network filters these external
353 fluctuations.

354 In contrast, the two dimensional spatial structure in our model supports rich chaotic network
355 dynamics outside the asynchronous state, yet with population-wide correlations that are a reasonable
356 mimic of experiment (Fig. 2aiv,c). Spatiotemporal chaos is a hallmark feature of systems that are far
357 from equilibrium in physics, chemistry and biology⁴¹. In particular, low viscosity fluids produce a
358 special brand of spatiotemporal chaotic behavior labelled turbulence, characterized by the presence
359 of vortices and eddies in the fluid flow⁵⁰. Like our network, the character of turbulent flow is very

360 dependent upon the dimension of the fluid, with one dimensional fluids not showing turbulence, and
361 two dimensional turbulent flow having larger spatial scales than the flow in full three dimensional
362 fluids⁵⁰. The dynamics within recurrent networks of neurons are certainly not equivalent to that of
363 fluids, in part because they possess both short and long range interactions in contrast to only the direct
364 local interactions in fluids. Nevertheless, the fluid analogy to our work is tempting since the chaotic
365 dynamics of our two dimensional network has a macroscopic character that permits low, but non-
366 vanishing, microscopic correlations, in contrast to the unrealistic high correlation dynamics of one
367 dimensional or disordered networks. While the top-down attentional signal in our model is similar
368 to that used in simpler models^{7,19}, the effect of top-down attention is to not only shift the operating
369 point of the network but also dampen the macroscopic chaotic dynamics of the network. In other
370 words, attention not only attenuates the transfer of population-wide variability but also quenches the
371 variability that is to be transferred. This permits a near complete attention-mediated suppression of
372 internally generated correlations (Fig. 6). This extreme sensitivity allows top-down inputs to easily
373 control the processing state of a network.

374 State dependent shifts in population-wide variability are widespread throughout cortex³, and are
375 often a signature of cognitive control. The circuit structure of our network is not a special feature of
376 the primate visual system, yet rather a generic property of most cortices. We thus expect that the basic
377 mechanisms for population-wide variability and its modulation exposed in our study will be operative
378 in many regions of the cortex, and in many animal systems.

379 **Methods**

380 **Network model description** The network consists of three layers. Layer 1 is modeled by a pop-
381 ulation of $N_1 = 2,500$ excitatory neurons, the spikes of which are taken as independent Poisson
382 processes with a uniform rate $r_1 = 10$ Hz. Layer 2 and Layer 3 are recurrently coupled networks
383 with excitatory ($\alpha = e$) and inhibitory ($\alpha = i$) populations of $N_e = 40,000$ and $N_i = 10,000$ neu-

384 rons, respectively. Each neuron is modeled as an exponential integrate-and-fire (EIF) neuron whose
385 membrane potential is described by:

$$C_m \frac{dV_j^\alpha}{dt} = -g_L (V_j^\alpha - E_L) + g_L \Delta_T e^{(V_j^\alpha - V_T)/\Delta_T} + I_j^\alpha(t). \quad (1)$$

386 Each time $V_j^\alpha(t)$ exceeds a threshold V_{th} , the neuron spikes and the membrane potential is held for
387 a refractory period τ_{ref} then reset to a fixed value V_{re} . Neuron parameters for excitatory neurons are
388 $\tau_m = C_m/g_L = 15$ ms, $E_L = -60$ mV, $V_T = -50$ mV, $V_{th} = -10$ mV, $\Delta_T = 2$ mV, $V_{re} = -65$ mV
389 and $\tau_{ref} = 1.5$ ms. Inhibitory neurons are the same except $\tau_m = 10$ ms, $\Delta_T = 0.5$ mV and $\tau_{ref} = 0.5$
390 ms. The total current to each neuron is:

$$\frac{I_j^\alpha(t)}{C_m} = \sum_{k=1}^{N_F} \frac{J_{jk}^{\alpha F}}{\sqrt{N}} \sum_n \eta_F(t - t_n^{F,k}) + \sum_{\beta=e,i} \sum_{k=1}^{N_\beta} \frac{J_{jk}^{\alpha \beta}}{\sqrt{N}} \sum_n \eta_\beta(t - t_n^{\beta,k}) + \mu_\alpha, \quad (2)$$

391 where $N = N_e + N_i$ is the total number of the network population. Postsynaptic current is

$$\eta_\beta(t) = \frac{1}{\tau_{\beta d} - \tau_{\beta r}} \begin{cases} e^{-t/\tau_{\beta d}} - e^{-t/\tau_{\beta r}}, & t \geq 0 \\ 0, & t < 0 \end{cases} \quad (3)$$

392 where $\tau_{er} = 1$ ms, $\tau_{ed} = 5$ ms and $\tau_{ir} = 1$ ms, $\tau_{id} = 8$ ms. The feedforward synapses from Layer
393 1 to Layer 2 have the same kinetics as the recurrent excitatory synapse, i.e. $\eta_F^{(2)}(t) = \eta_e(t)$. The
394 feedforward synapses from Layer 2 to Layer 3 have a fast and a slow component.

$$\eta_F^{(3)}(t) = p_f \eta_e(t) + p_s \eta_s(t)$$

395 with $p_f = 0.2$, $p_s = 0.8$. $\eta_s(t)$ has the same form as Eq. 3 with a rise time constant $\tau_r^s = 2$ ms and
396 a decay time constant $\tau_d^s = 100$ ms. The excitatory and inhibitory neurons in Layer 3 receive static
397 current μ_e and μ_i , respectively.

398 Neurons on the three layers are arranged on a uniform grid covering a unit square $\Gamma = [0, 1] \times [0, 1]$.
399 The probability that two neurons, with coordinates $\mathbf{x} = (x_1, x_2)$ and $\mathbf{y} = (y_1, y_2)$ respectively, are
400 connected depends on their distance measured periodically on Γ :

$$p_{\alpha\beta}(\mathbf{x}, \mathbf{y}) = \bar{p}_{\alpha\beta} g(x_1 - y_1; \alpha_\beta) g(x_2 - y_2; \alpha_\beta). \quad (4)$$

401 Here $\bar{p}_{\alpha\beta}$ is the mean connection probability and

$$g(x; \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \sum_{k=-\infty}^{\infty} e^{-(x+k)^2/(2\sigma^2)} \quad (5)$$

402 is a wrapped Gaussian distribution. Excitatory and inhibitory recurrent connection widths of Layer
 403 2 are $\alpha_{\text{rec}}^{(2)} := \alpha_e^{(2)} = \alpha_i^{(2)} = 0.1$ and feedforward connection width from Layer 1 to Layer 2 is
 404 $\alpha_{\text{ffwd}}^{(2)} = 0.05$. The recurrent connection width of Layer 3 is $\alpha_{\text{rec}}^{(3)} = 0.2$ and the feedforward connection
 405 width from Layer 2 to Layer 3 is $\alpha_{\text{ffwd}}^{(3)} = 0.1$. A presynaptic neuron is allowed to make more than
 406 one synaptic connection to a single postsynaptic neuron.

407 The recurrent connectivity of Layer 2 and Layer 3 have the same synaptic strengths and mean
 408 connection probabilities. The recurrent synaptic weights are $J_{ee} = 80$ mV, $J_{ei} = -240$ mV, $J_{ie} = 40$
 409 mV and $J_{ii} = -300$ mV. Recall that individual synapses are scaled with $1/\sqrt{N}$ (Eq. 2); so that,
 410 for instance, $J_{ee}/\sqrt{N} \approx 0.36$ mV. The mean connection probabilities are $\bar{p}_{ee} = 0.01$, $\bar{p}_{ei} = 0.04$,
 411 $\bar{p}_{ie} = 0.03$, $\bar{p}_{ii} = 0.04$. The out-degrees are $K_{ee}^{\text{out}} = 400$, $K_{ei}^{\text{out}} = 1600$, $K_{ie}^{\text{out}} = 300$ and $K_{ii}^{\text{out}} = 400$.
 412 The feedforward connection strengths from Layer 1 to Layer 2 are $J_{eF}^{(2)} = 140$ mV and $J_{iF}^{(2)} = 100$
 413 mV with probabilities $\bar{p}_{eF}^{(2)} = 0.1$ and $\bar{p}_{iF}^{(2)} = 0.05$ (out-degrees $K_{eF2}^{\text{out}} = 4000$ and $K_{iF2}^{\text{out}} = 500$). The
 414 feedforward connection strengths from Layer 2 to Layer 3 are $J_{eF}^3 = 25$ mV and $J_{iF}^3 = 15$ mV with
 415 mean probabilities $\bar{p}_{eF}^{(3)} = 0.05$ and $\bar{p}_{iF}^{(3)} = 0.05$ (out-degrees are $K_{eF3}^{\text{out}} = 2000$ and $K_{iF3}^{\text{out}} = 500$). Only
 416 the excitatory neurons in Layer 2 project to Layer 3.

417 The spatial models in Fig. 2a,ii, a,iv contain only Layer 1 and Layer 2. In the model with disordered
 418 connectivity, the connection probability between a pair of neurons is $\bar{p}_{\alpha\beta}$, independent of distance.
 419 Other parameters are the same as the spatial model. The decay time constant of IPSC (τ_{id}) was varied
 420 from 1 to 15 ms (Fig. 2d). The rise time constant of IPSC (τ_{ir}) is 1 ms when $\tau_{\text{id}} > 1$ ms and 0.5 ms
 421 when $\tau_{\text{id}} = 1$ ms.

422 The parameters used in Fig. 3c,d are $\mu_i = [0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4]$ pA and $\mu_E = 0$ pA.
 423 The mean firing rates in Layer 2 are $r_e^{(2)} = 19$ Hz and $r_i^{(2)} = 9$ Hz. In the further analysis (Fig. 4b-d,

424 Fig. 5b and Fig. 6), we used $\mu_I = 0.2$ pA for the unattended state and $\mu_I = 0.35$ pA for the attended
 425 state. In simulations of the spatial model with fast inhibition (Fig. 4c), $\tau_{ir} = 0.5$ ms, $\tau_{id} = 1$ ms. In
 426 simulations of the spatial model with broad inhibitory projection (Fig. 4d and Fig. 5c), $\alpha_e^{(3)} = 0.1$,
 427 $\alpha_i^{(3)} = 0.2$. Other parameters are not changed.

428 All simulations were performed on the CNBC Cluster in the University of Pittsburgh. All simula-
 429 tions were written in a combination of C and Matlab (Matlab R 2015a, Mathworks). The differential
 430 equations of the neuron model were solved using forward Euler method with time step 0.01 ms.

431 **Neural field model and stability analysis** We use a two dimensional neural field model to describe
 432 the dynamics of population rate (Fig. 5). The neural field equations are

$$\tau_\alpha \frac{\partial r_\alpha(x, t)}{\partial t} = -r_\alpha + \phi(w_{\alpha e} * r_e + w_{\alpha i} * r_i + \mu_\alpha) \quad (6)$$

433 where $r_\alpha(x, t)$ is the firing rate of neurons in population $\alpha = e, i$ near spatial coordinates $x \in [0, 1] \times$
 434 $[0, 1]$. The symbol $*$ denotes convolution in space, μ_α is a constant external input and $w_{\alpha\beta}(x) =$
 435 $\bar{w}_{\alpha\beta} g(x; \sigma_\beta)$ where $g(x; \sigma_\beta)$ is a two-dimensional wrapped Gaussian with width parameter σ_β , $\beta =$
 436 e, i . The transfer function is a threshold-quadratic function, $\phi(x) = [x^2]_+$. The timescale of synaptic
 437 and firing rate responses are implicitly combined into τ_α . In networks with approximate excitatory-
 438 inhibitory balance, rates closely track synaptic currents¹⁸, so τ_α represents the synaptic time constant
 439 of population $\alpha = e, i$.

440 For constant inputs, μ_e and μ_i , there exists a spatially uniform fixed point, which was computed
 441 numerically using an iterative scheme²⁵. Linearizing around this fixed point in Fourier domain gives
 442 a Jacobian matrix at each spatial Fourier mode²⁵

$$J(\vec{n}) = \begin{bmatrix} (-1 + g_e \tilde{w}_{ee}(\vec{n})) / \tau_e & g_e \tilde{w}_{ei}(\vec{n}) / \tau_e \\ g_i \tilde{w}_{ie}(\vec{n}) / \tau_i & (-1 + g_i \tilde{w}_{ii}(\vec{n})) / \tau_i \end{bmatrix}.$$

443 where $\vec{n} = (n_1, n_1)$ is the two-dimensional Fourier mode, $\tilde{w}_{\alpha\beta}(\vec{n}) = \bar{w}_{\alpha\beta} \exp(-2\|\vec{n}\|^2 \pi^2 \sigma_\beta^2)$ is the
 444 Fourier coefficient of $w_{\alpha\beta}(x)$ with $\|\vec{n}\|^2 = n_1^2 + n_2^2$ and g_a is the gain, which is equal to $\phi'(r_\alpha)$ evaluated

445 at the fixed point. The fixed point is stable at Fourier mode \vec{n} if both eigenvalues of $J(\vec{n})$ have negative
446 real part. Note that stability only depends on the wave number, $k = \|\vec{n}\|$, so Turing-Hopf instabilities
447 always occur simultaneously at all Fourier modes with the same wave number (spatial frequency).

448 For the stability analysis in Fig. 5a, τ_i varies from 2.5 ms to 25 ms, σ_i varies from 0.05 to 0.2,
449 and $\tau_e = 5$ ms and $\sigma_e = 0.1$. The rest of the parameters were $\bar{w}_{ee} = 80$, $\bar{w}_{ei} = -160$, $\bar{w}_{ie} = 120$,
450 $\bar{w}_{ii} = -200$, $\mu_e = 0.48$ and $\mu_i = 0.32$. Depolarizing the inhibitory population ($\mu_I = 0.5$) expands
451 the stable region (Fig. 5a, black dashed).

452 **Experimental methods** Each of the two datasets (recordings from V4 and recordings from V1
453 and MT) was collected from two different rhesus monkeys as they performed an orientation-change
454 detection task. All animal procedures were in accordance with the Institutional Animal Care and Use
455 Committee of Harvard Medical School, University of Pittsburgh and Carnegie Mellon University.

456 For analysis in Fig. 1b and Fig. 4a, data was collected with two microelectrode arrays implanted
457 bilaterally in area V4²⁴. In our analysis, we include stimulus presentations prior to the change stimulus
458 from correct trials, excluding the first stimulus in a trial to avoid adaptation effects. Spike counts
459 during the sustained response (120 - 260 ms after stimulus onset) are considered for the correlation
460 and factor analysis. Neurons recorded from either the left or right hemisphere in one session are
461 treated separately. There are a total of 42,496 trials for 72,765 pairs from 74 recording sessions. Two
462 sessions from the original study were excluded for factor analysis due to inadequate trials. The trial
463 number and unit number of each session is summarized in Table S1.

464 For analysis in Fig. 1d, data was collected with one microelectrode array implanted in area V1 and
465 a single electrode or a 24-channel linear probe inserted into MT²³. Again, our analysis includes full
466 contrast stimulus presentations prior to the change stimulus from correct trials and excludes the first
467 stimulus in a trial to avoid adaptation effects. Spike counts are measured 30 - 230 ms after stimulus
468 onset for V1 and 50 - 250 ms after stimulus onset for MT to account for the average visual latencies

469 of neurons in both areas. There are a total of 1,631 V1-MT pairs from 32 recording sessions.

470 **Statistical methods** To compute the noise correlation of each simulation, 500 neurons were ran-
471 domly sampled without replacement from the excitatory population of Layer 3 and Layer 2 within
472 a $[0, 0.5] \times [0, 0.5]$ square (considering periodic boundary condition). Spike counts were computed
473 using a sliding window of 200 ms with 1 ms step size and the Pearson correlation coefficients were
474 computed between all pairs. Neurons of firing rates less than 2 Hz were excluded from the compu-
475 tation of correlations. In Fig. 3c,d, for each μ_i there were 50 simulations and each simulation was
476 20 sec long. Connectivity matrices and the initial states of each neuron's membrane potential were
477 randomized in each simulation. The first 1 second of each simulation was excluded from the correla-
478 tion analysis. Standard error was computed based on the mean correlations of each simulation. For
479 simulations of Fig. 2d, there was one simulation of 20 seconds per τ_{id} and the connectivity matrices
480 were randomized for each simulation. To compute the noise correlation, 1000 neurons were randomly
481 sampled without replacement in the excitatory population of Layer 2 within a $[0, 0.5] \times [0, 0.5]$ square.
482 Correlations are computed between firing rates that are smoothed with a Gaussian window of width
483 10 ms.

484 Factor analysis assumes spike counts of n simultaneously recorded neurons $x \in \mathbb{R}^{n \times 1}$ is a multi-
485 variable Gaussian process

$$x \sim \mathcal{N}(\mu, LL^T + \Psi)$$

486 where $\mu \in \mathbb{R}^{n \times 1}$ is the mean spike counts, $L \in \mathbb{R}^{n \times m}$ is the loading matrix of the m latent variables
487 and $\Psi \in \mathbb{R}^{n \times 1}$ is a diagonal matrix of independent variances for each neuron. We choose $m = 5$ and
488 compute the eigenvalues of LL^T , λ_i ($i = 1, 2, \dots, 5$), ranked in descending order. We compute the
489 residual covariance after subtracting the first mode as

$$Q = \text{Cov}(x, x) - L_1 \times L_1'$$

490 where $\text{Cov}(x, x)$ is the raw covariance matrix of x and L_1 is the loading matrix when fitting with

491 $m = 1$. The mean raw covariance and residual (Fig. 4a-d, bottom) are the mean of the off-diagonal
492 elements of $\text{Cov}(x, x)$ and Q , respectively. When applying factor analysis on model simulations (Fig.
493 4b-d), we randomly selected 50 excitatory neurons from Layer 3, whose firing rates were larger than
494 2 Hz in both the unattended and attended states. There were 10 non-overlapping sampling of neurons
495 and we applied factor analysis on each sampling of neuron spike counts. There were 15 simulations
496 with fixed connectivity matrices, each of which was 20 seconds long. Spike trains were truncated
497 into 140 ms spike count window with a total of 2,025 counts per neuron. In simulations with fast
498 inhibition (Fig. 4c) and broad inhibitory projection (Fig. 4d), the feedforward connectivity from
499 Layer 2 to Layer 3 was the same as the one in simulations of the original model (Fig. 4b).

500 To study the chaotic population firing rate dynamics of Layer 3 (Fig. 6), we fixed the spike
501 trains realizations from Layer 1 neurons, the membrane potential states of the Layer 2 neurons and
502 all connectivity matrices. Only the initial membrane potentials of Layer 3 neurons were randomized
503 across trials. There were 10 realizations of Layer 1 and Layer 2, each of which was 20 sec long. For
504 each simulation of Layer 2, 20 repetitions with different initial conditions were simulated for Layer
505 3. The connectivity matrices in Layer 3 were the same across the 20 repetitions but different for each
506 realization of Layer 1 and Layer 2. The realizations of Layer 1 and Layer 2 and the connectivity
507 matrices were the same for the attended and unattended states. Trial-to-trial variance of Layer 3
508 population rates (Fig. 6d) was the variance of the mean population rates of the Layer 3 excitatory
509 population, smoothed by a 200 ms rectangular filter, across the 20 repetitions. The first second of
510 each simulation was discarded.

511 **Code availability** Computer code for all simulations and analysis of the resulting data is included
512 in Supplementary Software.

513 **Data availability** The data that support the findings of this study are available from the correspond-
514 ing author upon request.

515 **References**

- 516 1. Kass, R. E., Amari, S.-I., Arai, K., Brown, E. N., Diekman, C. O., Diesmann, M., Doiron, B.,
517 Eden, U. T., Fairhall, A. L., Fiddymment, G. M., et al. Computational neuroscience: Mathematical
518 and statistical perspectives. *Annual Review of Statistics and Its Application* **5**(1), 183–214 (2018).
- 519 2. Cohen, M. and Kohn, A. Measuring and interpreting neuronal correlations. *Nature neuroscience*
520 **14**(7), 811–819 (2011).
- 521 3. Doiron, B., Litwin-Kumar, A., Rosenbaum, R., Ocker, G., and Josic, K. The mechanics of state
522 dependent neural correlations. *Nature neuroscience* **19**(3), 383–393 (2016).
- 523 4. Lin, I.-C., Okun, M., Carandini, M., and Harris, K. D. The nature of shared cortical variability.
524 *Neuron* **87**(3), 644–656 (2015).
- 525 5. Rabinowitz, N. C., Goris, R. L., Cohen, M., and Simoncelli, E. Attention stabilizes the shared
526 gain of v4 populations. *eLife*, e08998 (2015).
- 527 6. Ecker, A. S., Berens, P., Cotton, R. J., Subramaniyan, M., Denfield, G. H., Cadwell, C. R.,
528 Smirnakis, S. M., Bethge, M., and Tolias, A. S. State dependence of noise correlations in macaque
529 primary visual cortex. *Neuron* **82**(1), 235–248 (2014).
- 530 7. Kanashiro, T., Ocker, G. K., Cohen, M. R., and Doiron, B. Attentional modulation of neuronal
531 variability in circuit models of cortex. *eLife* **6** (2017).
- 532 8. Williamson, R. C., Cowley, B. R., Litwin-Kumar, A., Doiron, B., Kohn, A., Smith, M. A., and
533 Byron, M. Y. Scaling properties of dimensionality reduction for neural populations and network
534 models. *PLOS Computational Biology* **12**(12), e1005141 (2016).
- 535 9. Scholvinck, M. L., Saleem, A. B., Benucci, A., Harris, K. D., and Carandini, M. Cortical state

- 536 determines global variability and correlations in visual cortex. *Journal of Neuroscience* **35**(1),
537 170–178 (2015).
- 538 10. Kelly, R. C., Smith, M. A., Kass, R. E., and Lee, T. S. Local field potentials indicate network
539 state and account for neuronal response variability. *Journal of computational neuroscience* **29**(3),
540 567–579 (2010).
- 541 11. Middleton, J. W., Omar, C., Doiron, B., and Simons, D. J. Neural correlation is stimulus modu-
542 lated by feedforward inhibitory circuitry. *The Journal of Neuroscience* **32**(2), 506–518 (2012).
- 543 12. Okun, M., Steinmetz, N. A., Cossell, L., Iacaruso, M. F., Ko, H., Barthó, P., Moore, T., Hofer,
544 S. B., Mrsic-Flogel, T. D., Carandini, M., et al. Diverse coupling of neurons to populations in
545 sensory cortex. *Nature* **521**(7553), 511–515 (2015).
- 546 13. Schmitz, T. W. and Duncan, J. Normalization and the cholinergic microcircuit: A unified basis
547 for attention. *Trends in Cognitive Sciences* , 1–16 (2018).
- 548 14. Ni, A., Ruff, D., Alberts, J., Symmonds, J., and Cohen, M. Learning and attention reveal a general
549 relationship between population activity and behavior. *Science* **359**(6374), 463–465 (2018).
- 550 15. Shadlen, M. N. and Newsome, W. T. The variable discharge of cortical neurons: implications for
551 connectivity, computation, and information coding. *The Journal of neuroscience* **18**(10), 3870–
552 3896 (1998).
- 553 16. van Vreeswijk, C. and Sompolinsky, H. Chaos in neuronal networks with balanced excitatory
554 and inhibitory activity. *Science* **274**, 1724–1726 (1996).
- 555 17. Amit, D. J. and Brunel, N. Model of global spontaneous activity and local structured activity
556 during delay periods in the cerebral cortex. *Cerebral cortex* **7**(3), 237–252 (1997).

- 557 18. Renart, A., de La Rocha, J., Bartho, P., Hollender, L., Parga, N., Reyes, A., and Harris, K. D. The
558 asynchronous state in cortical circuits. *Science* **327**(5965), 587–590 (2010).
- 559 19. Hennequin, G., Ahmadian, Y., Rubin, D. B., Lengyel, M., and Miller, K. D. Stabilized supra-
560 linear network dynamics account for stimulus-induced changes of noise variability in the cortex.
561 *bioRxiv* , 094334 (2016).
- 562 20. Ponce-Alvarez, A., Thiele, A., Albright, T. D., Stoner, G. R., and Deco, G. Stimulus-dependent
563 variability and noise correlations in cortical mt neurons. *Proceedings of the National Academy of*
564 *Sciences* **110**(32), 13162–13167 (2013).
- 565 21. Wimmer, K., Compte, A., Roxin, A., Peixoto, D., Renart, A., and De La Rocha, J. Sensory
566 integration dynamics in a hierarchical network explains choice probabilities in cortical area mt.
567 *Nature communications* **6** (2015).
- 568 22. Latham, P. E. Correlations demystified. *Nature neuroscience* **20**(1), 6 (2017).
- 569 23. Ruff, D. A. and Cohen, M. R. Attention increases spike count correlations between visual cortical
570 areas. *Journal of Neuroscience* **36**(28), 7523–7534 (2016).
- 571 24. Cohen, M. and Maunsell, J. Attention improves performance primarily by reducing interneuronal
572 correlations. *Nature neuroscience* **12**(12), 1594–1600 (2009).
- 573 25. Rosenbaum, R. and Doiron, B. Balanced networks of spiking neurons with spatially dependent
574 recurrent connections. *Physical Review X* **4**(2), 021039 (2014).
- 575 26. Rosenbaum, R., Smith, M., Kohn, A., Rubin, J., and Doiron, B. The spatial structure of correlated
576 neuronal variability. *Nature Neuroscience* **20**, 107–114 (2016).
- 577 27. Monteforte, M. and Wolf, F. Dynamic flux tubes form reservoirs of stability in neuronal circuits.
578 *Physical Review X* **2**(4), 041007 (2012).

- 579 28. Koch, C. *Biophysics of computation: information processing in single neurons*. Oxford university
580 press, (2004).
- 581 29. Pyle, R. and Rosenbaum, R. Spatiotemporal dynamics and reliable computations in recurrent
582 spiking neural networks. *Physical Review Letters* **118**(1), 018103 (2017).
- 583 30. Ecker, A. S., Denfield, G. H., Bethge, M., and Tolias, A. S. On the structure of neuronal popu-
584 lation activity under fluctuations in attentional state. *Journal of Neuroscience* **36**(5), 1775–1789
585 (2016).
- 586 31. Levy, R. B. and Reyes, A. D. Spatial profile of excitatory and inhibitory synaptic connectivity in
587 mouse primary auditory cortex. *Journal of Neuroscience* **32**(16), 5609–5619 (2012).
- 588 32. Horvat, S., Gamanut, R., Ercsey Ravasz, M., Magrou, L., Gamanut, B., Van Essen, D. C.,
589 Burkhalter, A., Knoblauch, K., Toroczkai, Z., and Kennedy, H. Spatial embedding and wiring
590 cost constrain the functional layout of the cortical network of rodents and primates. *PLoS biology*
591 **14**(7), e1002512 (2016).
- 592 33. Mariño, J., Schummers, J., Lyon, D. C., Schwabe, L., Beck, O., Wiesing, P., Obermayer, K., and
593 Sur, M. Invariant computations in local cortical networks with balanced excitation and inhibition.
594 *Nature neuroscience* **8**(2), 194 (2005).
- 595 34. Lim, S. and Goldman, M. S. Balanced cortical microcircuitry for spatial working memory based
596 on corrective feedback control. *Journal of Neuroscience* **34**(20), 6790–6806 (2014).
- 597 35. Kuchibhotla, K. V., Gill, J. V., Lindsay, G. W., Papadoyannis, E. S., Field, R. E., Sten, T. A. H.,
598 Miller, K. D., and Froemke, R. C. Parallel processing by cortical inhibition enables context-
599 dependent behavior. *Nature neuroscience* **20**(1), 62–71 (2017).

- 600 36. Kim, H., Ährlund-Richter, S., Wang, X., Deisseroth, K., and Carlén, M. Prefrontal parvalbumin
601 neurons in control of attention. *Cell* **164**(1), 208–218 (2016).
- 602 37. Cohen, M. and Maunsell, J. Using neuronal populations to study the mechanisms underlying
603 spatial and feature attention. *Neuron* **70**(6), 1192–1204 (2011).
- 604 38. Cunningham, J. P. and Byron, M. Y. Dimensionality reduction for large-scale neural recordings.
605 *Nature neuroscience* **17**(11), 1500–1509 (2014).
- 606 39. Mariño, J., Schummers, J., Lyon, D. C., Schwabe, L., Beck, O., Wiesing, P., Obermayer, K., and
607 Sur, M. Invariant computations in local cortical networks with balanced excitation and inhibition.
608 *Nature neuroscience* **8**(2), 194–201 (2005).
- 609 40. Ermentrout, B. Neural networks as spatio-temporal pattern-forming systems. *Reports on progress*
610 *in physics* **61**(4), 353 (1998).
- 611 41. Cross, M. C. and Hohenberg, P. C. Pattern formation outside of equilibrium. *Reviews of modern*
612 *physics* **65**(3), 851 (1993).
- 613 42. London, M., Roth, A., Beeren, L., Häusser, M., and Latham, P. E. Sensitivity to perturbations in
614 vivo implies high noise and suggests rate coding in cortex. *Nature* **466**(7302), 123 (2010).
- 615 43. Sato, T. K., Nauhaus, I., and Carandini, M. Traveling waves in visual cortex. *Neuron* **75**(2),
616 218–229 (2012).
- 617 44. Keane, A. and Gong, P. Propagating waves can explain irregular neural dynamics. *Journal of*
618 *Neuroscience* **35**(4), 1591–1605 (2015).
- 619 45. Compte, A., Sanchez-Vives, M. V., McCormick, D. A., and Wang, X.-J. Cellular and network
620 mechanisms of slow oscillatory activity (≈ 1 hz) and wave propagations in a cortical network
621 model. *Journal of neurophysiology* **89**(5), 2707–2725 (2003).

- 622 46. Bosking, W. H., Zhang, Y., Schofield, B., and Fitzpatrick, D. Orientation selectivity and the
623 arrangement of horizontal connections in tree shrew striate cortex. *Journal of neuroscience* **17**(6),
624 2112–2127 (1997).
- 625 47. Stringer, C., Pachitariu, M., Steinmetz, N. A., Okun, M., Bartho, P., Harris, K. D., Sahani, M.,
626 and Lesica, N. A. Inhibitory control of correlated intrinsic variability in cortical networks. *Elife*
627 **5**, e19695 (2016).
- 628 48. Bondy, A. G., Haefner, R. M., and Cumming, B. G. Feedback determines the structure of corre-
629 lated variability in primary visual cortex. *Nature neuroscience* **21**(3), 598–606 (2018).
- 630 49. Gómez-Laberge, C., Smolyanskaya, A., Nassi, J. J., Kreiman, G., and Born, R. T. Bottom-up and
631 top-down input augment the variability of cortical neurons. *Neuron* **91**(3), 540–547 (2016).
- 632 50. Davidson, P. *Turbulence: an introduction for scientists and engineers*. Oxford University Press,
633 USA, (2015).

634 **Supplemental Information**

635 Figures. S1-S4

636 Table. S1

637 Supplementary Methods

638 Captions for Movies S1-S4

639 **Acknowledgments**

640 NIH Grants CRCNS R01DC015139-01ZRG1 (B.D.), 4R00EY020844- 03 (M.R.C.), R01 EY022930
641 (M.R.C.), 5T32NS7391-14 (D.A.R.), and Core Grant P30 EY008098; NSF Grants DMS-1517828,

642 DMS-1654268 and DBI-1707400 (R.R.) and DMS-1517082 (B.D.); Vannevar Bush faculty fellow-
643 ship N00014-18-1-2002 (B.D.), a Whitehall Foundation Grant (M.R.C.); a Klingenstein-Simons Fel-
644 lowship (M.R.C.); grants from the Simons Foundation (B.D. and M.R.C.); a Sloan Research Fellow-
645 ship (M.R.C.); a McKnight Scholar Award (M.R.C.).

646 **Author Contributions**

647 C.H., M.R.C., and B.D. conceived the project; C.H. performed the simulations and data analysis; R.P.
648 and R.R analyzed the firing rate model; D.A.R. and M.R.C. provided the experimental data; B.D.
649 supervised the project; all authors contributed to writing the manuscript.

650 **Author Information**

651 The authors declare no competing financial interests. Correspondence and requests for materials
652 should be addressed to B.D. (bdoiron@pitt.edu).