PLOS SUBMISSION

The correlated state in balanced neuronal networks

Cody Baker¹, Christopher Ebsch¹, Ilan Lampl,² Robert Rosenbaum^{1,3*}

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

2 Department of Neurobiology, Weizmann Institute of Science, Rehovot, 7610001, Israel3 Interdisciplinary Center for Network Science and Applications, University of Notre Dame, Notre Dame, IN, USA

* Robert.Rosenbaum@nd.edu

Abstract

Understanding the magnitude and structure of inter-neuronal correlations and their relationship to synaptic connectivity structure is an important and difficult problem in computational neuroscience. Early studies show that neuronal network models with excitatory-inhibitory balance naturally create very weak spike train correlations. Later work showed that, under some connectivity structures, balanced networks can produce larger correlations between some neuron pairs, even when the average correlation is very small. All of these previous studies assume that the local neuronal network receives feedforward synaptic input from a population of uncorrelated spike trains. We show that when spike trains providing feedforward input are correlated, the downstream recurrent neuronal network produces much larger correlations. We provide an in-depth analysis of the resulting "correlated state" in balanced networks and show that, unlike the asynchronous state of previous work, it produces "tight" excitatory-inhibitory balance, consistent with in vivo cortical recordings.

Author summary

Correlation and synchrony between the activity of neurons in the brain is known to play a crucial role in the dynamics and coding properties of neuronal networks, and also mediates synaptic plasticity and learning. Therefore, it is important to understand the relationship between the structure of connectivity in a neuronal networks and the correlations between the activity of neurons in the network. Previous theoretical work shows that this relationship is constrained by the widely observed balance between excitatory (positive) and inhibitory (negative) input received by neurons in the network. We extend this previous theoretical work to account for the fact that inputs coming from outside the local neuronal network might come from neural populations that are themselves correlated or partially synchronous. Including this biologically realistic assumption changes the basic operating state of the network and produces a tighter balance between excitatory and inhibitory synaptic inputs that is consistent with in vivo recordings.

Introduction

Correlations between the spiking activity of cortical neurons have important consequences for neural dynamics and coding [1–3]. A quantitative understanding of

> how spike train correlations are generated and shaped by the connectivity structure of neural circuits is made difficult by the noisy and nonlinear dynamics of recurrent neuronal network models [4–7]. Linear response and related techniques have been developed to overcome some of these difficulties [8–14], but their accuracy typically require an assumption of sparse and/or weak connectivity and, in some models, an additional assumption that neurons receive uncorrelated, feedforward Gaussian white noise input. However, cortical circuits are densely connected and receive spatially and temporally correlated synaptic input from outside the local circuit [15–18].

> An alternative approach to analyzing correlated variability in recurrent neuronal network models is motivated in part by the widely observed balance between excitatory and inhibitory synaptic inputs in cortex [19-26]. When synaptic weights are scaled like $1/\sqrt{N}$ where N is the size of a model network, a cortex-like balance between excitation and inhibition arises naturally at large network size, which defines the "balanced state" [27, 28]. Early work on balanced networks assumed sparse connectivity to produce weak spike train correlations, but it was later shown that keeping connection probabilities $\mathcal{O}(1)$ naturally produces weak, $\mathcal{O}(1/N)$, spike train correlations, defining the "asynchronous state" [29]. While these extremely weak spike train correlations are consistent with some cortical recordings [30], the magnitude of correlations in cortex can depend on stimulus, cortical area, layer, and behavioral or cognitive state, and can be much larger than predicted by the asynchronous state [6, 31-35]. This raises the question of how larger correlation magnitudes can arise in balanced cortical circuits. Later theoretical work showed that larger correlations can be obtained between some cell pairs in densely connected networks with specially constructed connectivity structure [36–39], offering a potential explanation of the larger correlations often observed in recordings. These previous theoretical studies of correlated variability in balanced networks assume that the recurrent network receives feedforward synaptic input from an external population of uncorrelated spike trains, so feedforward input correlations arise solely from overlapping feedforward synaptic projections. In reality, feedforward synaptic input to a cortical population comes from thalamic projections, other cortical areas, or other cortical layers in which spike trains could be correlated.

> We extend the theory of densely connected balanced networks to account for correlations between the spike trains of neurons in an external, feedforward input layer. We show that correlations between the feedforward synaptic input to neurons in the recurrent network are $\mathcal{O}(N)$ in this model, but cancel with $\mathcal{O}(N)$ correlations between recurrent synaptic input to produce $\mathcal{O}(1)$ total input correlation and $\mathcal{O}(1)$ spike train correlations on average, defining what we refer to as the "correlated state" in densely connected balanced networks. This correlated state offers an alternative explanation for the presence of moderately large spike train correlations in cortical recordings. We derive a simple, closed form approximation for the average cross-spectral density between neurons' spike trains in the correlated state in term of synaptic parameters alone, without requiring the use of linear response theory or any other knowledge of neurons' transfer functions. We show that the tracking of excitatory synaptic input currents by inhibitory currents is more precise and more similar to in vivo recordings [22] in the correlated state than in the asynchronous state. We also investigate the applicability of linear response approximations to correlated variability in densely connected balanced networks. Taken together, our results extend the theory of correlated variability in balanced networks to the biologically realistic assumption that presynaptic neural populations are themselves correlated.

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

27

29

30

31

32

33

34

35

37

39

40

41

42

43

44

45

46

47

48

49

50

Results

SUBMISSION

We consider recurrent networks of N integrate-and-fire model neurons, N_e of which are excitatory and N_i inhibitory. Neurons are randomly and recurrently interconnected and also receive random feedforward synaptic input from an external population of N_x neurons whose spike trains are homogeneous Poisson processes with rate r_x (Fig. 1A).

The membrane potential of neuron j in population a = e, i obeys the exponential integrate-and-fire (EIF) dynamics

$$C_m \frac{dV_j^a}{dt} = -g_L(V_j^a - E_L) + D_T e^{(V_j^a - V_T)/D_T} + T_j^a(t)$$

with the added condition that each time $V_j^a(t)$ exceeds V_{th} , it is reset to V_{re} and a spike is recorded. We additionally set a lower bound on the membrane potential at $V_{lb} = -100$ mV. Spike trains are represented as a sum of Dirac delta functions,

$$S_j^a(t) = \sum_n \delta(t - t_n^{a,j}),$$

where $t_n^{a,j}$ is the *n*th spike time of neuron *j* in population a = e, i, x. The total synaptic input current to neuron *j* in population a = e, i is decomposed as

$$T_{j}^{a}(t) = E_{j}^{a}(t) + I_{j}^{a}(t) + X_{j}^{a}(t)$$

where

$$B_{j}^{a}(t) = \sum_{k=1}^{N_{b}} J_{jk}^{ab}(\alpha_{b} * S_{j}^{b})(t)$$
(1)

for B = E, I, X and b = e, i, x respectively where * denotes convolution, J_{jk}^{ab} is the synaptic weight from neuron k in population b to neuron j in population a, and $\alpha_b(t)$ is a postsynaptic current (PSC) waveform. Without loss of generality, we assume that $\int \alpha_b(t) = 1$. We use $\alpha_b(t) = \tau_b^{-1} e^{-t/\tau_b} H(t)$ where H(t) is the Heaviside step function, though our results do not depend sensitively on the precise neuron model or PSC kernel used. For calculations, it is useful to decompose the total synaptic input into its recurrent and external sources,

$$T_j^a(t) = R_j^a(t) + X_j^a(t)$$

where

$$R_i^a(t) = E_i^a(t) + I_i^a(t)$$

is the recurrent synaptic input from the local circuit.

Local cortical circuits contain a large number of neurons and individual cortical neurons receive synaptic input from thousands of other neurons within their local circuit and from other layers or areas. Densely connected balanced networks have been proposed to model such large and densely interconnected neuronal networks [29,38]. In such models, one considers the limit of large N (with N_x , N_e and N_i scaled proportionally) with fixed connection probabilities and where synaptic weights are scaled like $\mathcal{O}(1/\sqrt{N})$ [28,29]. This scaling naturally captures the balance of mean excitatory and mean inhibitory synaptic input, as well as the tracking of excitation by inhibition, observed in cortical recordings [29]. In particular, we consider a random connectivity structure in which

$$J_{jk}^{ab} = \frac{1}{\sqrt{N}} \begin{cases} j_{ab} & \text{with probability } p_{ab} \\ 0 & \text{otherwise} \end{cases}$$
(2)

PLOS

72

73

74

75

76

77

79

81

82

83

55 56 57

58

52

53

54

64

65

66

67

68

69

70

SUBMISSION

where connections are statistically independent and $j_{ab}, p_{ab} \sim \mathcal{O}(1)$ for b = e, i, x and a = e, i. We furthermore define the proportions 85

$$q_b = \frac{N_b}{N}$$

which are assumed $\mathcal{O}(1)$. For all examples we consider, $q_e = 0.8$ and $q_i = q_x = 0.2$.

We next introduce notational conventions for quantifying the statistics of spike trains and synaptic inputs in the network. The mean firing rates of neurons in population a = e, i, x is defined by r_a for a = e, i, x and it is useful to define the 2×1 vector, $\boldsymbol{r} = [r_e \ r_i]^T$. The mean is technically interpreted as the expectation over realizations of the network connectivity, but for large N it is approximately equal to the sample mean over all neurons the network. Similarly, mean-field synaptic inputs to neurons in populations a = e, i are defined by

$$\bar{U}_a = \mathrm{mean}_i U_i^a(t)$$

for U = E, I, X, R, T and, in vector form, $\overline{U} = [U_e \ U_i]^T$

For quantifying correlated variability, we use the cross-spectral density (CSD)

$$\langle U^a_j, Z^b_k \rangle(f) = \int_{-\infty}^\infty C_{U^a_j Z^b_k}(\tau) e^{-2\pi i f \tau} d\tau$$

between $U_i^a(t)$ and $Z_k^b(t)$ for U, Z = E, I, X, S, R, T and a, b = e, i, x where

$$C_{U_i^a, Z_h^b}(\tau) = \operatorname{cov}(U_i^a(t), Z_k^b(t+\tau))$$

is the cross-covariance function. The argument, f, is the frequency at which the CSD is 97 evaluated. The CSD is a convenient measure of correlated variability because it simplifies mathematical calculations due to the fact that it is a Hermitian operator and 99 because most commonly used measures of correlated variability can be written as a 100 function of the CSD. For example, the cross-covariance is the inverse Fourier transform 101 of the CSD. Spike count covariances over large time windows can be written in terms of 102 the CSD by first noting that the spike count is an integral of the spike train [4], 103

spike count over
$$[0, t_0] = \int_0^{t_0} S_j^a(t) dt.$$

For large t_0 , the cross-spectrum between two integrals is related to the zero-frequency 104 CSD, 105

$$\lim_{t_0 \to \infty} \frac{1}{t_0} \operatorname{cov}\left(\int_0^{t_0} U_j^a(t) dt, \int_0^{t_0} Z_k^b(t) dt\right) = \langle U_j^a, Z_k^b \rangle (f=0).$$
(3)

Hence,

spike count covariance over $[0, t_0] \approx t_0 \langle S_i^a, S_k^b \rangle (f = 0).$

Following this result, we often quantify covariability between spike trains and between 107 synaptic currents using the zero-frequency CSD, which we estimate by taking the 108 covariance between integrals as in Eq. (3) using $t_0 = 250$ ms. This provides a simple, 109 easily estimated quantity for quantifying covariance. 110

Most of our computations are performed at the level of population averages, so we define

$$U_a, W_b \rangle = \operatorname{mean}_{j \neq k} \langle U_j^a, W_k^b \rangle.$$

which is a scalar function of frequency, f, for each a, b = e, i, x and 113 U, W = E, I, X, S, R, T. It is also convenient to define the 2×2 mean-field matrix form, 114

$$\langle \boldsymbol{U}, \boldsymbol{W} \rangle = \begin{bmatrix} \langle U_e, W_e \rangle & \langle U_e, W_i \rangle \\ \langle U_i, W_e \rangle & \langle U_i, W_i \rangle \end{bmatrix}.$$

106

111

112

87

89

90

91

92

93

95



for U, W = E, I, X, S, R, T. We also define the recurrent and feedforward mean-field ¹¹⁵ connectivity matrices, ¹¹⁶

$$W = \left[\begin{array}{cc} w_{ee} & w_{ei} \\ w_{ie} & w_{ii} \end{array} \right] \text{ and } W_x = \left[\begin{array}{cc} w_{ex} \\ w_{ix} \end{array} \right]$$

where $w_{ab}(f) = p_{ab}j_{ab}q_b\widetilde{\eta}_b(f) \sim \mathcal{O}(1)$ with $\widetilde{\eta}_b(f)$ the Fourier transform of $\eta_b(t)$. For the exponential kernels we use, $\widetilde{\eta}_b(f) = 1/(1 + 2\pi i f \tau_b)$. The zero-frequency values, $\overline{w}_{ab} = w_{ab}(0) = p_{ab}j_{ab}q_b$, define time-averaged interactions and mean-field connection matrices, $\overline{W} = W(0)$ and $\overline{W}_x = W_x(0)$.

This choice of notation allows us to perform computations on mean-field spike train and input statistics in a mathematically compact way. To demonstrate this, we first review the mean-field analysis of firing rates in the balanced state [27, 28, 40–42]. Mean external input is given by $\overline{\mathbf{X}} = \sqrt{N} \ \overline{W}_x r_x$ and mean recurrent input by $\overline{\mathbf{R}} = \sqrt{N} \ \overline{W} \mathbf{r}$ so that mean total synaptic input is given by

$$\overline{\boldsymbol{T}} = \sqrt{N} \left[\overline{W} \boldsymbol{r} + \overline{W}_x r_x \right].$$

In the balanced state, \overline{T} , $r \sim \mathcal{O}(1)$, which can only be obtained by a cancellation between external and recurrent synaptic inputs. This cancellation requires $\overline{W}r + \overline{W}_x r_x \sim \mathcal{O}(1/\sqrt{N})$ so that [27, 28, 40-42] 128

$$\lim_{N \to \infty} \boldsymbol{r} = -\overline{W}^{-1} \overline{W}_x r_x \tag{4}$$

in the balanced state. Hence, the balanced state can only be realized when this solution has positive entries, $r_e, r_i > 0$, which requires that [27, 28, 40]

 $\overline{X}_e/\overline{X}_i > w_{ei}/w_{ii} > w_{ee}/w_{ie}$. Below, we perform analogous derivations of mean-field CSDs in balanced networks.

A review of the asynchronous balanced state

We first review prior theoretical work that derives mean-field CSDs when spike trains in the external population are uncorrelated Poisson processes (Fig. 1A,B), so

$$\langle S_x, S_x \rangle = 0.$$

Since the derivations of these results [38] and analogous derivations for networks of binary neuron models [29] are presented elsewhere and since the derivations are similar to those presented below, we only review the results here and give the details of the derivation in Materials and Methods.

Since spike trains in the external population are uncorrelated, correlations between the external input to neurons in the recurrent network arise solely from overlapping feedforward synaptic projections with [29,38] (see Materials and Methods)

$$\langle \boldsymbol{X}, \boldsymbol{X} \rangle = q_x^{-1} W_x r_x W_x^* \sim \mathcal{O}(1).$$
⁽⁵⁾

where W_x^* is the conjugate transpose of W_x . If the spike trains in the external layer were uncorrelated, but not Poisson, then r_x in this equation could be replaced by the power spectral density of the spike trains.

It would at first seem that this $\mathcal{O}(1)$ external input correlation would lead to $\mathcal{O}(1)$ correlations between neurons' spike trains, but this is prevented by a cancellation between positive and negative sources of input correlation. In particular, correlations between neurons' recurrent synaptic inputs, $\langle \boldsymbol{R}, \boldsymbol{R} \rangle$, are also positive and $\mathcal{O}(1)$, but these positive sources of input correlations are canceled by negative correlations

133

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

117

118

119

Fig 1. The asynchronous state in densely connected balanced networks. A) Network diagram. An external population, X, of uncorrelated Poisson processes provides feedforward input to a randomly connected recurrent network of excitatory, E, and inhibitory, I, neurons. Feedforward input correlations are solely from overlapping projections from X. B,C) Raster plot of 200 randomly selected neurons from population X and E respectively. **D)** Histogram of external (X, green) recurrent (R = E + I, purple) and total (T = X + E + I, black) to all excitatory neurons. Currents here and elsewhere are reported in units $C_m V/s$ where C_m is the arbitrary membrane capacitance. E) Mean external (green), recurrent (purple), and total (black) input to excitatory neurons for networks of different sizes, N. F) Mean excitatory (red) and inhibitory (blue) neuron firing rates for difference network sizes. Solid curves are from simulations and dashed curves are from Eq. (4). G) Mean covariance between pairs of excitatory neurons' external inputs (green), recurrent inputs (purple), total inputs (black), and mean covariance between the recurrent input to one excitatory neuron and external input to the other (yellow) for different network sizes. Covariances were computed by integrating the inputs over 250ms windows then computing covariances between the integrals, which is proportional to zero-frequency CSD and has a closer relationship with spike count covariance (see Eq. (3) and surrounding discussion). Integrated currents have units $C_m mV$, so their covariances have units $C_m^2 m V^2$. H) Zoomed in view of black curve from E on a log-log axis (mean total input covariance, black) plotted alongside the function c/N (dashed gray) where c was chosen so that the two curves match at the largest N value. I) Mean spike count covariance between excitatory neuron spike trains (red), between inhibitory neuron spike trains (blue), and between excitatory-inhibitory pairs of spike trains (purple). Counts were computed over 250ms time windows. Solid curves are from simulations, dashed from Eq. (7) evaluated at zero frequency. Network size was $N = 10^5$ in B-D.

between neurons' recurrent and external inputs, $\langle X, R \rangle$, in such a way that the total synaptic input correlation is weak, 151

$$\langle \boldsymbol{T}, \boldsymbol{T} \rangle = \langle \boldsymbol{X}, \boldsymbol{X} \rangle + \langle \boldsymbol{X}, \boldsymbol{R} \rangle + \langle \boldsymbol{R}, \boldsymbol{X} \rangle + \langle \boldsymbol{R}, \boldsymbol{R} \rangle \sim \mathcal{O}(1/N)$$

where $\langle \mathbf{R}, \mathbf{X} \rangle = \langle \mathbf{X}, \mathbf{R} \rangle^*$. This cancellation is realized when the mean-field CSD between spike trains satisfies [38]

$$\langle \boldsymbol{S}, \boldsymbol{S} \rangle = \frac{1}{N} W^{-1} \langle \boldsymbol{X}, \boldsymbol{X} \rangle W^{-*} + \frac{1}{N} C_0 + o(1/N)$$
(6)

where $N \times o(1/N) \to 0$ as $N \to \infty$ and W^{-*} is the inverse of W^* . The first term in 155 this equation represents spike train correlations inherited from external inputs, namely 156 $\langle \boldsymbol{X}, \boldsymbol{X} \rangle$. The second term, C_0 , represents correlations generated intrinsically by chaotic 157 or chaos-like dynamics in the network [28, 38, 40, 43]. For example, a network with 158 deterministic, constant external input, $X_i^a(t) = \overline{X}_a$, generates correlated variability [40] 159 despite the fact that such networks are deterministic once an initial condition is 160 specified and therefore $\langle \mathbf{X}, \mathbf{X} \rangle = 0$ for such networks. While an exact expression for C_0 161 is unknown, we show empirically below that its effects are small compared to 162 correlations inherited from external input, at least for the network parameters that we 163 consider. Therefore, the following approximation is relatively accurate 164

$$\langle \boldsymbol{S}, \boldsymbol{S} \rangle \approx \frac{1}{N} W^{-1} \langle \boldsymbol{X}, \boldsymbol{X} \rangle W^{-*}.$$
 (7)

To demonstrate these results, we first simulated a network of $N = 10^4$ randomly and recurrently connected neurons receiving feedforward input from a population of

153

> Fig 2. The correlated state in densely connected balanced networks. A,B) Same as Fig. 1 B,C except spike trains in the external population, X, were correlated Poisson processes with spike count correlation c = 0.1. C) Mean CSD between excitatory neuron spike trains (red), between inhibitory neuron spike trains (blue), and between excitatory-inhibitory pairs (purple). Solid curves are from simulations (each CSD averaged over 10^5 pairs) and dashed are from Eq. (14). D-I) Same as Fig. 1D-I except spike trains in the external population, X, were correlated Poisson processes with spike count correlation c = 0.1. Dashed lines in I are from Eq. (14) evaluated at zero frequency. Network size was $N = 10^5$ in A-D.

 $N_x = 2000$ uncorrelated Poisson-spiking neurons (Fig. 1A,B). As predicted, spiking activity in the recurrent network was asynchronous and irregular (Fig. 1C; mean spike count correlation between neurons with rates at least 1 Hz was 5.2×10^{-4}) with approximate balance between external (X) and recurrent (R) synaptic input sources (Fig. 1D). Varying the network size, N, demonstrates the $\mathcal{O}(\sqrt{N})$ growth of mean external (\overline{X}) and recurrent (\overline{R}) synaptic input currents that cancel to produce $\mathcal{O}(1)$ mean total input current (\overline{T}) (Fig. 1E), as predicted by the mean-field theory of balance. As a result, firing rates converge to the limiting values predicted by Eq. (4) (Fig. 1F).

As predicted by the analysis of the asynchronous state, the mean covariances between individual sources of synaptic inputs appear $\mathcal{O}(1)$ (Fig. 1G), but cancel to produce much smaller, $\mathcal{O}(1/N)$, total input covariance (Fig. 1G,H). Mean spike count covariances also appear $\mathcal{O}(1/N)$ and show good agreement with the closed form approximation from Eq. (7) (Fig. 1I).

The fact that the approximation in Eq. (7) accurately captures the scaling of spike count covariances from simulations implies that correlated variability inherited by overlapping synaptic inputs dominate intrinsic correlations, C_0 , which are ignored in Eq. (7). Similar results were found in previous work [38]. This contrasts with previous findings in networks of binary neurons, in which intrinsic variability appeared to dominate [44].

The correlated state in balanced networks

Above, we analyzed correlated variability when spike trains in the external population were uncorrelated, $\langle S_x, S_x \rangle = 0$, which produced asynchronous spiking in the recurrent network, $\langle S, S \rangle \sim \mathcal{O}(1/N)$. We now relax this assumption by considering moderate correlations between neurons in the external layer (Fig. 2A),

$$\langle S_x, S_x \rangle \sim \mathcal{O}(1).$$

Most previous work analyzing spike train correlations in recurrent networks relies on 191 knowledge of the "correlation susceptibility" or "transfer" function, which quantifies the 192 mapping from synaptic input covariance to spike train covariance, e.g., the mapping 193 from $\langle T_e, T_e \rangle$ to $\langle S_e, S_e \rangle$. However, susceptibility functions depend sensitively on the 194 neuron model being used and their derivation typically relies on diffusion 195 approximations that assume neurons' input currents can be approximated by Gaussian 196 white noise and are weakly correlated [8, 11]. These assumptions are not valid for the 197 densely connected, correlated networks with exponentially-decaying PSC kernels we 198 consider here. 199

Instead of assuming knowledge of a covariance transfer or susceptibility function, our derivation relies only on an assumption that transfer of mean-field covariance is $\mathcal{O}(1)$, specifically that 200

$$\frac{\langle S_a, U_b \rangle}{\langle T_a, U_b \rangle} \sim \mathcal{O}(1). \tag{8}$$

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

SUBMISSION

In other words, mean-field spiking statistics are not drastically different in magnitude

<

than mean-field input statistics because neurons transfer inputs to outputs in an $\mathcal{O}(1)$ fashion. Importantly, we do not need to know the value of the fraction in Eq. (8).

In addition to this assumption, all of our derivations follow from a few simple arithmetical rules that rely on the bilinearity of the operator $\langle \cdot, \cdot \rangle$. Specifically (see Materials and Methods for derivations),

$$\langle \boldsymbol{T}, \boldsymbol{U} \rangle = \langle \boldsymbol{R}, \boldsymbol{U} \rangle + \langle \boldsymbol{X}, \boldsymbol{U} \rangle$$
 (9)

203

204

205

206

207

208

209

210

211

212

215

216

217

218

219

220

221

224

225

$$\langle \boldsymbol{R}, \boldsymbol{U} \rangle = \sqrt{N} W \langle \boldsymbol{S}, \boldsymbol{U} \rangle \tag{10}$$

 $\langle \boldsymbol{X}, \boldsymbol{U} \rangle = \sqrt{N} W_x \langle S_x, \boldsymbol{U} \rangle \tag{11}$

$$\langle \boldsymbol{U}, \boldsymbol{Z} \rangle = \langle \boldsymbol{Z}, \boldsymbol{U} \rangle^*$$
 (12)

for any $U, Z = E, I, X, R, S, S_x, T$ where A^* is the conjugate-transpose of A and where we omit smaller order terms here and below (see Materials and Methods for details). Eq. (9) follows from the fact that total input is composed of recurrent and external sources, T = R + X. Eqs. (10) and (11) follow from the fact that recurrent and external inputs are composed of linear combinations of $\mathcal{O}(N)$ spike trains, *c.f.* Eq. (1), and that synaptic weights are $\mathcal{O}(1/\sqrt{N})$. Eq. (12) is simply a property of the Hermitian cross-spectral operator.

We first derive the CSD between external inputs to neurons in the recurrent network. ²¹³ Applying of Eqs. (11) and (12) gives ²¹⁴

$$\begin{aligned} \langle \boldsymbol{X}, \boldsymbol{X} \rangle &= \sqrt{NW_x} \langle S_x, \boldsymbol{X} \rangle \\ &= NW_x \langle S_x, S_x \rangle W_x^* \\ &\sim \mathcal{O}(N) \end{aligned}$$

Hence, $\mathcal{O}(1)$ covariance between the spike trains in the external population induces $\mathcal{O}(N)$ covariance between the external input currents to neurons in the recurrent network. This is a result of the effects of "pooling" on correlations and covariances, namely that the covariance between two sums of N correlated random variables is typically $\mathcal{O}(N)$ times larger than the covariances between the individual summed variables [29, 45, 46].

We next derive the CSD between spike trains and external inputs. First note that

$$T, X\rangle = \langle R, X \rangle + \langle X, X \rangle$$

= $\sqrt{N}W \langle S, X \rangle + \langle X, X \rangle.$ (13)

However, it follows from our assumption that neuronal transfer is $\mathcal{O}(1)$ (see Eq. (8)) that $\langle T, X \rangle \sim \langle S, X \rangle$. Therefore, we have 223

 $\langle \boldsymbol{S}, \boldsymbol{X} \rangle \sim \sqrt{N} W \langle \boldsymbol{S}, \boldsymbol{X} \rangle + \langle \boldsymbol{X}, \boldsymbol{X} \rangle,$

which is only consistent if there is a cancellation between the two terms on the right hand side. Specifically, we must have that

$$\langle \boldsymbol{S}, \boldsymbol{X} \rangle = -\frac{1}{\sqrt{N}} W^{-1} \langle \boldsymbol{X}, \boldsymbol{X} \rangle \sim \mathcal{O}(\sqrt{N})$$

$$egin{aligned} &\langle m{S},m{T}
angle = \sqrt{N}\langlem{S},m{S}
angle W^* + \langlem{S},m{X}
angle \ &= \sqrt{N}\langlem{S},m{S}
angle W^* - rac{1}{\sqrt{N}}W^{-1}\langlem{X},m{X}
angle. \end{aligned}$$



However, our assumption of $\mathcal{O}(1)$ transfer implies that $\langle \boldsymbol{S}, \boldsymbol{T} \rangle \sim \langle \boldsymbol{S}, \boldsymbol{S} \rangle$ so

$$\langle \boldsymbol{S}, \boldsymbol{T} \rangle \sim \sqrt{N} \langle \boldsymbol{S}, \boldsymbol{S} \rangle W^* - \frac{1}{\sqrt{N}} W^{-1} \langle \boldsymbol{X}, \boldsymbol{X} \rangle$$

which is only consistent if there is cancellation between the terms on the right hand side. ²³⁰ This cancellation can only be realized if ²³¹

$$\langle \boldsymbol{S}, \boldsymbol{S} \rangle = \frac{1}{N} W^{-1} \langle \boldsymbol{X}, \boldsymbol{X} \rangle W^{-*}$$

= $W^{-1} W_x \langle S_x, S_x \rangle W_x^* W^{-*}$ (14)

which is $\mathcal{O}(1)$ and where we have omitted terms of smaller order. Notably, evaluating Eq. (14) does not depend on knowledge of neurons' correlation susceptibility functions or any other neuronal transfer properties, but only depends on synaptic parameters and input statistics. 232

In summary, $\mathcal{O}(1)$ covariance between spike trains in the external population produces $\mathcal{O}(N)$ covariance between neurons' external inputs, but $\mathcal{O}(1)$ covariance between spike trains in the recurrent network on average. We hereafter refer to this state as the "correlated state" since it produces moderately strong spike train correlations in contrast to the asynchronous state characterized by extremely weak spike train correlations. The reduction from $\mathcal{O}(N)$ external input covariance to $\mathcal{O}(1)$ spike train covariance arises from a cancellation mechanism analogous to the one that reduces $\mathcal{O}(1)$ external input correlation to $\mathcal{O}(1/N)$ spike train correlations in the asynchronous state (see above).

To demonstrate these results, we simulated a network of $N = 10^4$ neurons identical to the network from Fig. 1 except that spike trains in the external population were correlated Poisson processes (Fig. 2A) with

$$\langle S_x, S_x \rangle(f) = cr_x e^{-4f^2 \pi^2 \tau_c^2}.$$
 (15)

Here, $r_x = 10$ Hz is the same firing rate used in Fig. 1 and c = 0.1 quantifies the spike count correlation coefficient between the spike trains in the external population over large counting windows. See Materials and Methods for a description of the algorithm used to generate the spike trains.

The recurrent network exhibited moderately correlated spike trains in contrast to spike trains in the asynchronous state (Fig. 2B, compare to Fig. 1C; mean spike count correlation between neurons with rates at least 1 Hz was 0.077). The mean CSDs between spike trains in the recurrent network closely matched the theoretical predictions from Eq. (14) (Fig. 2C). As in the asynchronous state, external and recurrent synaptic input sources approximately canceled (Fig. 2D), as predicted by balanced network theory.

Varying N demonstrates that the network exhibits the same cancellation between $\mathcal{O}(\sqrt{N})$ mean external and recurrent synaptic input sources and that Eq. (4) for the mean firing rates is accurate (Fig. 2E,F). As predicted by the analysis of the correlated state, the covariance between individual sources of input currents appear $\mathcal{O}(N)$ (Fig. 2G), but cancel to produce much smaller, approximately $\mathcal{O}(1)$, total input covariance (Fig. 2G,H). Mean spike count covariances also appear $\mathcal{O}(1)$ and converge toward the limit predicted by Eq. (14) (Fig. 2I). Hence, despite the complexity of spike timing dynamics in densely connected balanced networks, mean-field spike train covariances is accurately predicted by a simple, linear equation in terms of synaptic parameters. The derivation of this equation does not require the use of linear response theory, which can be problematic for densely connected networks with synaptic kinetics and non-vanishing correlations.

Fig 3. Excitatory-inhibitory tracking in vivo and in simulations. A) In vivo membrane potential recordings from neurons in rat barrel cortex, reproduced from [22]. Each pair of traces are simultaneously recorded membrane potentials. Red traces were recorded in current clamp mode near the reversal potential of inhibition and blue traces near the reversal potential of excitation (with action potentials pharmacologically suppressed), so red traces are approximately proportional to excitatory input current fluctuations and blue traces approximate inhibitory input current fluctuations. Vertical scale bars are 20mV. B,C) Excitatory (red) and inhibitory (blue) synaptic input currents to two randomly selected excitatory neurons in the asynchronous (B) and correlated (C) states. Simulations were the same as those in Figs. 1B-D and 2A-D respectively. Currents are plotted with outward polarity negative and inward positive.

The correlated state produces tight balance between excitatory and inhibitory input fluctuations consistent with cortical recordings

We have so far considered cancellation between positive and negative sources of input correlations at the mean-field level, *i.e.*, averaged over pairs of postsynaptic neurons (Figs. 1G,H and 2G,H). In vivo cortical recordings reveal that this cancellation occurs even at the level of single postsynaptic neuron pairs. When one neuron was clamped near its inhibitory reversal potential and another neuron is clamped near its excitatory reversal potential (with spiking suppressed), recorded membrane potential fluctuations are approximately mirror images of one another (Fig. 3A, top). Similarly, if both neurons are held near their excitatory reversal potential (Fig. 3A, bottom), recorded membrane potential fluctuations are their inhibitory reversal potential (Fig. 3A, bottom), recorded membrane potential fluctuations are highly correlated. This implies that fluctuations in the excitatory and inhibitory input to other nearby neurons (see [22] for more details and interpretation).

To test whether this phenomenon occurred in our simulations, we randomly chose two neurons and decomposed their synaptic input into the total excitatory (E + X) and the inhibitory (I) components. In the asynchronous state, input current fluctuations were fast and largely unshared between neurons or between current sources in the same neuron (Fig. 3B), in contrast to evidence from *in vivo* recordings. Input current fluctuations in the correlated state were slower, larger, and most importantly largely synchronized between neurons (Fig. 3C), consistent with evidence from *in vivo* recordings. This precise tracking of fluctuations in excitatory and inhibitory synaptic currents is referred to as "tight balance" [47] (as opposed to the "loose balance" of the asynchronous state). The results would be similar if we decomposed inputs into their external (X) and recurrent (R = E + I) sources instead of excitatory (E + X) and inhibitory (I).

To better understand this striking difference between input currents in the asynchronous and correlated states, we first computed the average covariance between the excitatory and inhibitory input to pairs of (excitatory) neurons in the network. These averages have the same dependence on network size, N, as they do when input currents are broken into external and recurrent sources (compare Fig. 4A,B to Figs. 1G and 2G). Specifically, in the asynchronous state, covariances between individual current sources are $\mathcal{O}(1)$ on average, but cancel to produce weak $\mathcal{O}(1/N)$ covariance between the total synaptic input to neurons on average (Fig. 4A). In the correlated state, the average covariance between individual input sources is $\mathcal{O}(N)$ and cancellation produces $\mathcal{O}(1)$ average total input covariance (Fig. 4B).

Hence, despite the precise cancellation of positive and negative sources of input

> Fig 4. The scaling of mean and variance of excitatory and inhibitory input covariance in the asynchronous and correlated states. A,B) Same as Figs. 1G and 2G, except inputs were decomposed into their excitatory (E + X), and inhibitory (I) components instead of external and recurrent. Red curves show mean excitatory-excitatory input covariance, blue show inhibitory-inhibitory, purple show excitatory-inhibitory, and black curves show total (same as black curves in Figs. 1G and 2G). C,D) Histogram of input current covariances across all excitatory cell pairs for a network of size $N = 10^5$. E,F) Same as A,B except we plotted the variance of covariances across cell pairs instead of the mean. As above, integrated currents have units $C_m mV$, so input covariances have units $C_m^2 mV^2$ and the variance of covariances have units $C_m^4 mV^4$ where C_m is the arbitrary membrane capacitance.

covariance at the mean-field level in the asynchronous state (Fig. 4A), this cancellation 310 is apparently not observed at the level of individual neuron pairs (Fig. 3B). To see why 311 this is the case, we computed the distribution of input current covariances across all 312 pairs of excitatory neurons. We found that these distributions were broad and the 313 distribution of total input covariance was highly overlapping with the distributions of 314 individual input current sources (Fig. 4C, the black distribution overlaps with the 315 others). This implies that cancellation does not reliably occur at the level of individual 316 pairs since, for example, the total input covariance for a pair of neurons can be similar in 317 magnitude or even larger than the covariance between those neurons' excitatory inputs. 318

The distributions of input covariances were strikingly difference in the correlated state. The distribution of total input covariances was far narrower than the distributions of individual input current sources and the distributions were virtually non-overlapping (Fig. 4D). Hence, a precise cancellation between positive and negative sources of input covariance must occur for every neuron pair, leading to the tight balance observed in Fig. 3C.

These results are more precisely understood by computing the variance across neuron pairs of input covariances as N is varied. In the asynchronous state, the variance of input covariances from all sources appear to scale like $\mathcal{O}(1)$ (Fig. 4E). Since the mean input covariance between individual sources are also $\mathcal{O}(1)$ (Fig. 4A), the overlap between distributions in Fig. 4C is expected. In the correlated state, the variances of input covariances appear to scale like $\mathcal{O}(N)$ except for the variance of the total input covariance, which appears to scale like $\mathcal{O}(1)$ (Fig. 4F). Since the variances scale like $\mathcal{O}(N)$, the standard deviations scale like $\mathcal{O}(\sqrt{N})$. This, combined with the fact that the mean input covariances between individual sources scale like $\mathcal{O}(N)$, implies that the distributions in Fig. 4E will be non-overlapping when N is large. The same conclusions would be reached if we decomposed inputs into their external (X) and recurrent (R = E + I) sources instead of total excitatory (X + E) and inhibitory (I).

Testing linear response approximations to pairwise spike train covariance

As discussed above, Eqs. (7) and (14) are appealing because, unlike many previous approximations of spike train CSD and spike count covariance in recurrent networks [8–11,13], they do not require knowledge of neurons' correlation susceptibility functions or any other neuronal transfer properties. However, they are limited because they only give the population-averaged covariance or CSD.

Previous studies that provide approximations for the full $N \times N$ matrix of all spike train CSDs or spike count covariances use linear response approximations that rely on assumptions that the network is sparsely or weakly coupled and/or that external input is uncorrelated or weakly correlated Gaussian white noise [10, 11]. These assumptions 346

PLOS

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

> permit the application of Fokker-Planck and linear response techniques. External input 348 in our models is not weakly correlated and is not approximated well by Gaussian white 349 noise, so these approaches are not fully justified in out model, but might nevertheless 350 vield a useful approximation. We next derive and test a linear response approximation 351 to spike train covariability in our model. 352

> First define $\vec{S}(t) = [S_1^e(t) \dots S_{N_e}^e(t) S_1^i(t) \dots S_{N_i}^i(t)]^T$ to be the full $N \times 1$ vector 353 of spike trains in the recurrent network obtained by concatenating the excitatory and 354 inhibitory spike train vectors. Define the $N \times 1$ synaptic input vectors $\vec{T}(t)$, $\vec{R}(t)$, and 355 X(t) similarly. The total input vector can be written as 356

$$\vec{T} = \vec{R} + \vec{X}$$

 $\vec{R} = J * \vec{S}$

with

SUBMISSION

and

$$\vec{X} = J_x * \vec{S}_x.$$

Here,

$$= \left[\begin{array}{cc} J_{ee} & J_{ei} \\ J_{ie} & J_{ii} \end{array} \right],$$

is the $N \times N$ recurrent connectivity matrix and

 $J_x = \left[\begin{array}{c} J_{ex} \\ J_{ie} \end{array} \right],$

is the $N \times N_x$ feedforward connectivity matrix with $N_a \times N_b$ blocks defined by 361 $[J_{ab}]_{jk}(t) = J^{ab}_{jk}\eta_b(t)$. Note that the connection weights are time dependent and * denotes the matrix product with multiplication replaced by convolution [11]. 362 363

The linearity of transfer from \vec{S} to \vec{R} and from \vec{S}_x to \vec{X} implies that

J

$$\langle \vec{R}, \vec{U} \rangle = \widetilde{J} \langle \vec{S}, \vec{U} \rangle \tag{16}$$

and

$$\langle \vec{X}, \vec{U}
angle = \widetilde{J}_x \langle \vec{S}_x, \vec{U}
angle$$

for $\vec{U} = \vec{S}, \vec{S}_x, \vec{R}, \vec{X}, \vec{T}$ where $\widetilde{J}(f)$ is the element-wise Fourier transform of the matrix 366 J(t) and similarly for $J_x(f)$. This implies, for example, that the $N \times N$ matrix of external input CSDs is given by

$$\langle \vec{X}, \vec{X} \rangle = \tilde{J}_x \langle \vec{S}_x, \vec{S}_x \rangle \tilde{J}_x^*$$

and similarly for the CSDs between recurrent inputs.

 $\langle \vec{R}, \vec{R} \rangle = \widetilde{J} \langle \vec{S}, \vec{S} \rangle \widetilde{J}^*$

In summary, the transfer of spike train CSDs to input CSDs follows easily from the 370 linearity of transfer from spike trains to inputs [4]. However, a closed expression for the 371 matrix of spike train CSDs in the recurrent network, $\langle \vec{S}, \vec{S} \rangle$, requires knowledge of the 372 transfer from total input to spike train CSDs. 373

If spike trains were related linearly to total input, $\vec{S} = A * \vec{T}$ for some diagonal matrix, A(t), then we could derive an exact closed equation for $\langle \vec{S}, \vec{S} \rangle$. However, integrate-and-fire neuron models transfer their input currents to spike trains nonlinearly.

368

357

350

360

364

365

369

374

375

SUBMISSION

Fig 5. Testing the accuracy of linear response approximations in densely connected balanced networks in the asynchronous and correlated states. **A,B)** Spike count covariance from simulations of a network with $N = 10^5$ in the asynchronous (A) and correlated (B) states plotted against the theoretical value given by evaluating Eq. (18) evaluated at zero frequency, with empirically estimated gains used for $\tilde{A}_j(0)$. **C,D)** The left-hand-side (LHS) of Eq. (19) plotted against the right-hand-side (RHS) using $\langle \vec{S}, \vec{S} \rangle$ estimated from the same simulations as in A,B and using the same empirically estimated gains. **E,F)** Same as C,D except we used the mean gain for all values of $\tilde{A}_j(0)$.

Nevertheless, we can obtain an approximation to spike train CSDs by assuming an approximate linear transfer of CSDs, specifically that [11]

$$\langle \vec{S}, \vec{X} \rangle \approx \widetilde{A} \langle \vec{T}, \vec{X} \rangle, \langle \vec{S}, \vec{T} \rangle \approx \widetilde{A} \langle \vec{T}, \vec{T} \rangle,$$

$$\langle \vec{S}, \vec{S} \rangle \approx \widetilde{A} \langle \vec{T}, \vec{S} \rangle$$

$$(17)$$

where $\widetilde{A}(f)$ is a diagonal matrix with $\widetilde{A}_{jj}(f)$ the "susceptibility function" of neuron j [8, 48-50]. Note that these equations would be exactly true if neural transfer were linear, $\vec{S} = A * \vec{T}$. It can be derived from these assumptions that (see Materials and Methods for derivation and [11] for similar derivations) 379

$$\langle \vec{S}, \vec{S} \rangle \approx [\widetilde{A}^{-1} - \widetilde{J}]^{-1} \langle \vec{X}, \vec{X} \rangle [\widetilde{A}^{-1} - \widetilde{J}]^{-*}.$$
(18)

One problem with testing the approximation in Eq. (18) is that we do not have an estimate of $\widetilde{A}_{jj}(f)$. Spike count covariances over large windows are proportional to zero-frequency CSD (see Eq. (3) and surrounding discussion). Evaluated at zero frequency, a neuron's susceptibility function is equal to its gain, *i.e.* the derivative of the neuron's firing rate with respect to its mean input [49, 50], 387

$$\widetilde{A}_j(0) = \frac{dr_j}{d\overline{T}_j}.$$

We therefore tested Eq. (18) for spike count covariances by estimating the gain empirically from simulations. In doing so, we found that Eq. (18) is only moderately accurate at approximating spike count covariances from simulations, both in the asynchronous state (Fig. 5A) and in the correlated state (Fig. 5B). We suspected that some of the error was due to numerical inaccuracies introduced by inverting the large, ill-conditioned matrices in Eq. (18). To test this hypothesis, we re-wrote Eq. (18) in a mathematically equivalent formulation that does not involve matrix inverses, 399

$$(Id - \widetilde{A}\widetilde{J})\langle \vec{S}, \vec{S}\rangle(Id - \widetilde{A}\widetilde{J}) \approx \widetilde{A}\langle \vec{X}, \vec{X}\rangle \widetilde{A}^*$$
(19)

where Id is the $N \times N$ identity matrix. We tested the accuracy of Eq. (19) at zero frequency using the same empirically estimated values of the gains and using the matrix, $\langle \vec{S}, \vec{S} \rangle (0)$, of spike count covariances estimated from simulations and found that it was very accurate, especially in the correlated state (Fig. 5C,D). Since Eqs. (18) and (19) are mathematically equivalent, this suggests that much of the observed inaccuracy of Eq. (18) is due to the presence of inverses of large, ill-conditioned matrices.

One shortcoming of Eqs. (18) and (19) is that they require an estimate of \hat{A} . We obtained this estimate empirically by simulating the entire network and numerically estimating the gains, which greatly limits the utility of the equations. We next

395

396

397

398

399

400

401

402

403

377

> wondered whether the accuracy of the equations is sensitive to the accuracy of the gain 404 estimate, or if a rough approximation to the gains would be sufficient. To answer this 405 question, we tested Eq. (19) again, but replaced the individual estimated gains on the 406 diagonal of A with their mean. In other words, we used $A(0) = \overline{q}Id$ where \overline{q} is the 407 average gain estimated from simulations and Id is the $N \times N$ identity matrix. This 408 replacement greatly decreased the accuracy of Eq. (19) (Fig. 5E,F), suggesting that an 409 accurate estimate of the individual gains is necessary for applying and interpreting 410 Eqs. (18) and (19). 411

In summary, linear response approximations are fairly accurate in the densely connected balanced networks with spatially and temporally correlated noisy feedforward input studied here (Fig. 5C,D). However, their utility is limited by two factors: First, that using linear response theory to approximate spike train CSDs or spike count covariances requires the inversion of large, ill-conditioned matrices, which introduces substantial error (Fig. 5A,B). Secondly, that the application of linear response approximation requires estimates of neurons' individual susceptibility functions or gains. The mean-field equations (7) and (14) do not have these shortcomings, but only give the population-averaged CSDs and covariances. It should also be noted that correlations were weak in our simulations, even in the correlated state (mean spike count correlations between neurons with firing rate above 1 Hz was 0.077). Stronger correlations can be obtained with alternative parameters, which could potentially make Eqs. (18) and (19) less accurate.

Correlated variability from singular mean-field connectivity structure

We have shown that $\mathcal{O}(1)$ spike train correlations can be obtained in balanced networks by including correlations between neurons in an external layer ($\langle S_x, S_x \rangle \sim \mathcal{O}(1)$), defining what we refer to as the "correlated state." Previous work shows that $\mathcal{O}(1)$ spike train correlations can be obtained in the recurrent network with uncorrelated external spike trains ($\langle S_x, S_x \rangle = 0$) when the mean-field connectivity matrix is constructed in such a way that the recurrent network cannot achieve the cancellation required for these states to be realized [37–39]. This is most easily achieved using networks with several discrete sub-populations or networks with distance-dependent connectivity. For simplicity, we restrict our analysis to discrete sub-populations. We first extend the mean-field theory from above to such networks, then generalize and extend the analysis from previous work to understand the emergence of $\mathcal{O}(1)$ correlations when the mean-field connectivity matrix is singular.

The recurrent networks considered above have two statistically homogeneous sub-populations: one excitatory and one inhibitory and the external population is a single homogeneous population. Suppose instead that there are K sub-populations in the recurrent network, with the kth population containing $N_k = q_k N$ neurons where $\sum_k q_k = 1$. Connectivity is random with p_{jk} denoting the connection probability from population k to j, and j_{jk}/\sqrt{N} denoting the strengths of the connections for $j, k = 1, \ldots, K$. All neurons in population k are assumed to have the PSC kernel $\eta_k(t)$ which is again assumed to have integral 1. Similarly, suppose that the external network contains K_x sub-populations each with $N_k^x = q_k^x N_x$ neurons where $q_k^x = \sum_k q_{x,k} = 1$. Feedforward connection probabilities and strengths are given by p_{jk}^x and j_{jk}^x/\sqrt{N} for $j = 1, \ldots, K$ and $k = 1, \ldots, K_x$. Assume that $q_k, p_{jk}, j_{jk}, q_k^x, p_{jk}^x$, and j_{jk}^x are all $\mathcal{O}(1)$. We then define the $K \times K$ mean-field recurrent connectivity matrix by $[W]_{jk} = p_{jk}j_{jk}q_k\tilde{\eta}_k$ and the mean-field feedforward connectivity matrix by $[W_x]_{jk} = p_{jk}^x j_{jk}^x q_k^x \tilde{\eta}_k^x$. For all of the networks considered above, we had K = 2 and $K_x = 1$.

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

When W is an invertible matrix, Eqs. (4), (7), and (14) are equally valid for networks with several subpopulations as they are for the simpler networks considered above. Hence, the mean-field theory of firing rates and correlations extends naturally to networks with several populations [38–42, 51]. However, when W is singular, Eqs. (4), (7), and (14) cannot be evaluated. Instead, Eq. (4) can be re-written as

$$W\boldsymbol{r} = -W_x \boldsymbol{r}_x. \tag{20}$$

459

460

461

462

463

464

465

466

467

471

472

473

474

475

476

477

478

479

480

When W is singular, this equation only has a solution, \mathbf{r} , when $\overline{\mathbf{X}} = -W_x \mathbf{r}_x$ is in the range or "column space" of W. Otherwise, balance is broken. An in-depth analysis of firing rates in such networks is provided in previous work [41,42,51] (and extended to spatially continuous networks in [40,51]), so we hereafter assume that $\overline{\mathbf{X}}$ is in the range of W and balance is achieved.

A similar analysis may be applied to spike train CSDs. For simplicity, we assume here that spike trains in the external population are uncorrelated, $\langle S_x, S_x \rangle = 0$, since this is the case considered in previous work and since this is the case in which a singular W breaks the asynchronous state. Eq. (7) can be re-written as

$$W\langle \boldsymbol{S}, \boldsymbol{S} \rangle W^* = \frac{1}{N} \langle \boldsymbol{X}, \boldsymbol{X} \rangle.$$
 (21)

where we have ignored smaller order terms. When W is singular, Eq. (21) is not guaranteed to have a solution, $\langle \boldsymbol{S}, \boldsymbol{S} \rangle$. More precisely, a solution can only exist when the $K \times K$ matrix, $\langle \boldsymbol{X}, \boldsymbol{X} \rangle$, is in the range of the linear operator \mathcal{L} defined by

$$\mathcal{L}U = WUW^*.$$

In that case, Eq. (21) has a solution so that $\langle \boldsymbol{S}, \boldsymbol{S} \rangle \sim \mathcal{O}(1/N)$ and the asynchronous state is still realized. However, if $\langle \boldsymbol{X}, \boldsymbol{X} \rangle$ is not in the range of \mathcal{L} , the asynchronous state cannot be realized because Eq. (21) does not have a solution.

Darshan et al. [39] looked at a similar scenario except the singularity of their networks made them incapable of cancelling *internally generated covariance*, in contrast to the *external input covariance* considered here. Other work [37,38] analyzed the scenario with external input covariance and singular connectivity, as well as the extension to spatially extended networks. However, this previous work did not completely analyze the structure of correlations, but only showed that the asynchronous state was broken. We next show that correlations in the recurrent network are $\mathcal{O}(1)$ and derive their structure.

Using Eqs. (9), (10), and (12), we can write the mean-field total input CSD as

$$\langle \boldsymbol{T}, \boldsymbol{T} \rangle = NW \langle \boldsymbol{S}, \boldsymbol{S} \rangle W^* + \sqrt{N} (W \langle \boldsymbol{S}, \boldsymbol{X} \rangle + \langle \boldsymbol{X}, \boldsymbol{S} \rangle W^*) + \langle \boldsymbol{X}, \boldsymbol{X} \rangle.$$
(22)

If W is not invertible, then W^* has a non-trivial nullspace. Let v_1, v_2, \ldots, v_n be a basis for the nullspace of W^* and define

$$P = v_1 v_1^* + v_2 v_1^* + \dots + v_n v_n^*$$

which is a self-adjoint matrix that defines the orthogonal projection onto the nullspace of W^* . Note that P is a Hermitian matrix $(P = P^*)$ and $PW = W^*P = 0$ (the zero matrix). Define the projection $A_0 = PAP$ for any matrix A. Unless $\langle X, X \rangle$ is carefully constructed otherwise, we can expect that

$$\langle \boldsymbol{X}, \boldsymbol{X} \rangle_0 \sim \langle \boldsymbol{X}, \boldsymbol{X} \rangle \sim \mathcal{O}(1).$$

Then take the projection of both sides of Eq. (22) above to get

$$\langle \boldsymbol{T}, \boldsymbol{T} \rangle_0 = \langle \boldsymbol{X}, \boldsymbol{X} \rangle_0 \sim \mathcal{O}(1)$$
 (23)

Fig 6. Correlated variability in a balanced network with singular

mean-field connectivity matrix. A) Network schematic. The recurrent network is statistically identical to the networks considered previously, but there are two external populations that each connect to a different half of the neurons in the recurrent network. B) Same as Fig. 1F, but for the multi-population network from A. C) Same as Fig. 1G, but for the network in A and where input covariances are only averaged over postsynaptic neurons in the same group (both postsynaptic cells in e_1 or both in e_2). The dashed gray curve shows the theoretical prediction for total input covariance (the black curve) from Eq. (24). D) Same as Fig. 1I, but for the network in A and where spike count covariances are only averaged over postsynaptic neurons in the same group (first cell in a_j and second cell in b_j for a, b = e, i and j = 1, 2). E,F) Same as C and D, but covariances are computed between cells in opposite groups (one cell in a_1 and the other cell in b_2).

where we have omitted terms of order $\langle \boldsymbol{X}, \boldsymbol{X} \rangle / N$ (see Materials and Methods for more details). Hence, the total input CSD is $\mathcal{O}(1)$ when $\langle \boldsymbol{X}, \boldsymbol{X} \rangle$ is not in the range of \mathcal{L} , even though it is $\langle \boldsymbol{X}, \boldsymbol{X} \rangle / N$ when W is invertible (*i.e.*, in the asynchronous state). Moreover, the structure of $\langle \boldsymbol{T}, \boldsymbol{T} \rangle$ is given to highest order in N by $\langle \boldsymbol{X}, \boldsymbol{X} \rangle_0 = P \langle \boldsymbol{X}, \boldsymbol{X} \rangle P$, which can be computed exactly from knowledge of the structure of $\langle \boldsymbol{X}, \boldsymbol{X} \rangle$ and W.

When neural transfer from T to S is $\mathcal{O}(1)$ (see Eq. (8) and surrounding discussion), this implies that $\langle S, S \rangle \sim \mathcal{O}(1)$ so that the asynchronous state is broken when $\langle X, X \rangle$ is not in the range of \mathcal{L} . While we cannot be certain that $\langle S, S \rangle$ has the same structure as $\langle T, T \rangle$, it should have a similar structure as long as neural transfer of correlations is similar for each sub-population.

To demonstrate these results, we consider the same network from above with re-wired feedforward projections from the external population. Specifically, divide the excitatory, inhibitory, and external populations each into two equal-sized sub-populations, labeled e_1 , i_1 , x_1 , e_2 , i_2 , and x_2 where population a_k contains $N_a/2$ neurons. Hence the network has the same total number of neurons as before, but we have simply sub-divided the populations. To distinguish this network from the one considered in Figs. 1 and 2, we refer to the previous network as the 3-population network and to this modified network as the 6-population network.

We re-wire the feedforward connections so that x_1 only connects to e_1 and i_1 , and x_2 only projects to e_2 and i_2 . Specifically, we set the connection probabilities to $p_{a_jx_k} = 2p_{ax}$ if j = k and $p_{a_jx_k} = 0$ if $j \neq k$ for a, b = e, i and j, k = 1, 2, where p_{ab} are the connection probabilities for the 3-population network and $p_{a_jb_k}$ for the 6-population network. This re-wiring assures that neurons in the recurrent network receive the same number of feedforward connections on average from the external population. The recurrent connectivity structure is not changed at all. Specifically, we set $p_{a_jb_k} = p_{ab}$ for a, b = e, i. All connection strengths are unchanged, $j_{a_jb_k} = j_{ab}$ for a = e, i and b = e, i, xand all neurons in the external population have the same firing rate, r_x , as before. See Fig. 6A for a schematic of this network.

The feedforward mean-field connectivity matrix can be written in block form as

$$W_x = \begin{bmatrix} W_x^{2 \times 1} & \mathbf{0} \\ \mathbf{0} & W_x^{2 \times 1} \end{bmatrix}$$

where **0** is the 2 × 1 zero-matrix and $W_x^{2\times 1} = [w_{ex} \ w_{ix}]^T$ is the 2 × 1 feedforward connectivity matrix for the 3-population network. Note that W_x is 4 × 2 since there are 4 populations in the recurrent network and 2 populations in the external population. The recurrent mean-field connectivity matrix is

$$W = \frac{1}{2} \left[\begin{array}{cc} W^{2\times 2} & W^{2\times 2} \\ W^{2\times 2} & W^{2\times 2} \end{array} \right]$$

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507



where

$$W^{2\times 2} = \left[\begin{array}{cc} w_{ee} & w_{ei} \\ w_{ie} & w_{ii} \end{array} \right]$$

is the 2 × 2 recurrent connectivity matrix for the 3-population network. Note that W is 509 4 × 4. Here, $w_{ab} = p_{ab}j_{ab}q_b\tilde{\eta}_b$ are the same values used above for analyzing the 510 3-population network. 511

Even though W is non-invertible, $\overline{X} = W_x [r_x \ r_x]^T$ is in the range of W for this example, so firing rates in the balanced state can be computed using Eq. (20), and are identical to the firing rates for the 3-population networks considered above.

[1]

The nullspace of W^* is spanned by the orthonormal vectors

and

$$v_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 0\\ -1\\ 0 \end{bmatrix}$$
$$v_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 0\\ 1\\ 0\\ -1 \end{bmatrix}$$

so the projection matrix is given in block form by

$$P = \frac{1}{2} \left[\begin{array}{cc} I_2 & -I_2 \\ -I_2 & I_2 \end{array} \right]$$

where I_2 is the 2×2 identity matrix.

The external input CSD is determined by the average number of overlapping feedforward projections to any pair of neurons in the recurrent network (multiplied by their connection strength and r_x), which gives (in block form)

$$\langle \boldsymbol{X}, \boldsymbol{X} \rangle = 2 \begin{bmatrix} \langle \boldsymbol{X}, \boldsymbol{X} \rangle^{2 \times 2} & \boldsymbol{0} \\ \boldsymbol{0} & \langle \boldsymbol{X}, \boldsymbol{X} \rangle^{2 \times 2} \end{bmatrix}$$

where **0** is the 2 × 2 zero matrix and $\langle \boldsymbol{X}, \boldsymbol{X} \rangle^{2 \times 2}$ is the external input CSD from the 3-population network, given by Eq. (5). Therefore, by Eq. (23),

$$\langle \boldsymbol{T}, \boldsymbol{T} \rangle_0 = \langle \boldsymbol{X}, \boldsymbol{X} \rangle_0 = P \langle \boldsymbol{X}, \boldsymbol{X} \rangle P = \begin{bmatrix} \langle \boldsymbol{X}, \boldsymbol{X} \rangle^{2 \times 2} & -\langle \boldsymbol{X}, \boldsymbol{X} \rangle^{2 \times 2} \\ -\langle \boldsymbol{X}, \boldsymbol{X} \rangle^{2 \times 2} & \langle \boldsymbol{X}, \boldsymbol{X} \rangle^{2 \times 2} \end{bmatrix}.$$
(24)

In other words, the mean total input CSD between excitatory neurons in the same 518 subgroup (two neurons in e_1 or two neurons in e_2 ; diagonal blocks above) is positive 519 and equal to half the mean external input between the same neurons. Hence, the 520 cancellation by the recurrent network only reduces the external input CSD by a factor 521 of 1/2, as opposed to the $\mathcal{O}(1/N)$ reduction realized in the asynchronous state (when W 522 is invertible). In contrast, the mean total input CSD between excitatory neurons in 523 opposite subgroups (one neuron in e_1 and the other in e_2 ; off-diagonal blocks above) has 524 the same magnitude as for same-subgroup pairs, but is negative. This represents a 525 competitive dynamic between the two groups since they inhibit one another (recurrent 526 connections are net-inhibitory in balanced networks [27, 42]), but receive different 527 feedforward input noise. Interestingly, the average CSD between all pairs of spike trains 528 is still $\mathcal{O}(1/N)$ in this example, but it is easy to design examples with singular W in 529 which this is not true. A similar example was considered in previous work [38], but 530 external input was generated artificially instead of coming from an external population. 531

515

516

Fig 7. Synaptic input currents in a balanced network with correlations from singular mean-field connectivity. Same as Fig. 3A except for the network from Fig. 6. The left two traces are input currents to two excitatory neurons in population e_1 (cells 1 and 2). The right traces are input currents to an excitatory neuron in population e_2 (cell 3).

Simulating this network for varying values of N shows that firing rates approach 532 those predicted by the balance equation (20) (Fig. 6B), confirming that balance is 533 realized. Pairs of excitatory neurons in the same group (both neurons in e_1 or both 534 neurons in e_2) receive positively correlated external input and recurrent input (Fig. 6C, 535 purple and green curves) that are partially canceled by negative correlations between 536 their recurrent and excitatory input (Fig. 6C, yellow curve). Because the cancellation is 537 only partial, the correlation between the neurons' total inputs is $\mathcal{O}(1)$ (Fig. 6C, black 538 curve) in contrast to the asynchronous state (compare to Fig. 1G,H where cancellation 539 is perfect at large N). The total input covariance agrees well with the theoretical 540 prediction from Eq. (24) (Fig. 6C, dashed gray line). As a result of this lack of 541 cancellation between total input covariance, spike count covariances are also $\mathcal{O}(1)$ and 542 positive between same-group pairs (Fig. 6D). For opposite group pairs (one neuron in e_1 543 and the other in e_2), cancellation is also imperfect, but this leads to negative total input 544 covariance, in agreement with the theoretical prediction from Eq. (24) (Fig. 6E), and 545 leads to negative spike count covariances between neurons in opposite populations 546 (Fig. 6F). 547

In summary, we have analyzed two mechanisms to generate $\mathcal{O}(1)$ spike train 548 correlations in balanced networks. For the first mechanism (Fig. 2), spike trains in the 549 external population are correlated so that external input correlations are $\mathcal{O}(N)$. 550 Cancellation is achieved so that spike train correlations are reduced to $\mathcal{O}(1)$. For the 551 other mechanism (Fig. 6), external input correlation is $\mathcal{O}(1)$, but precise cancellation 552 cannot be achieved so that spike trains inherit the $\mathcal{O}(1)$ correlations from the input. 553 How could these two mechanisms be distinguished in cortical recordings? Under the first 554 mechanism, we showed that fluctuations of inhibitory input to individual neurons closely 555 tracks fluctuations of other neurons' excitatory inputs (Fig. 3C). This should not be the 556 case under the second mechanism because precise cancellation is not realized. Indeed, 557 plotting the excitatory and inhibitory input to three excitatory neurons (two in e_1 and 558 one in e_2) shows that input fluctuations are not closely tracked (Fig. 7). This provides a 559 way to distinguish the two mechanisms from paired intracellular recordings. Indeed, the 560 first mechanism (which we refer to as the "correlated state") appears more consistent 561 with the cortical recordings considered here (compare Fig. 3A to Figs. 3C and 7). 562

Discussion

We analyzed correlated variability in recurrent, balanced networks of integrate-and-fire neurons receiving correlated feedforward input from an external population. We showed that correlations between spike trains in the recurrent network are small $(\mathcal{O}(1/N))$ when spike trains in the external population are uncorrelated, consistent with previous work on the asynchronous state [29,38], but much larger $(\mathcal{O}(1))$ when spike trains in the external population are correlated, giving rise to a "correlated state." In both states, strong correlations in the feedforward input are canceled by recurrent synaptic input due to the excitatory-inhibitory tracking that arises naturally in densely connected balanced networks. This cancellation allows for the derivation of a concise and accurate closed form expression for spike train CSDs in terms of synaptic parameters alone. Hence correlations in balanced networks are determined predominately by synaptic

563

564

565

566

567

568

570

571

572

573

574

LOS

SUBMISSION

PLOS SUBMISSION

connectivity structure, not neuronal dynamics. The tracking of excitatory synaptic input by inhibition was observable on a pair-by-pair basis in the correlated state, but not the asynchronous state, suggesting that the correlated state is more consistent with in vivo recordings.

We only considered recurrent networks with two, statistically homogeneous neural populations: one excitatory and one inhibitory. Our analysis can be extended to multiple subpopulations as long as each sub-population contains $\mathcal{O}(N)$ neurons, and also extends to networks with connection probabilities that depend on distance, orientation tuning, or other continuous quantities. This analysis has been developed for the asynchronous state in previous work [38] and is easily extended to the correlated state as well. The primary difference is that $\langle \mathbf{X}, \mathbf{X} \rangle$ is $\mathcal{O}(N)$ instead of $\mathcal{O}(1)$.

Previous work has shown that networks with multiple sub-populations and networks with distance-dependent connectivity can break the asynchronous state in balanced networks when the network connectivity structure is constructed in such a way that the recurrent network cannot achieve the cancellation required for the asynchronous state [37–39], leading to $\mathcal{O}(1)$ correlations between some cell pairs. We showed that the precise tracking of excitation by inhibition provides an experimentally testable prediction for distinguishing this mechanism from the one underlying the correlated state (see Fig. 7 and surrounding discussion).

Another alternative mechanism for achieving larger correlations in balanced networks is through instabilities of the balanced state. Such instabilities, especially in spatially extended networks, can create pattern-forming dynamics that produce correlated spiking without hypersynchrony [40,52–57]. Future studies should work towards experimentally testable predictions that distinguish correlations that arise from instabilities from those that arise through the mechanisms considered here. For example, since instabilities generate correlations internally, they should produce weak correlations between activity in the recurrent network and activity in the external population(s) providing input to that network [57], in contrast to the mechanisms we consider here. Indeed, some recordings show that local circuit connectivity can increase correlations [58], which is consistent with internally generated correlations, but inconsistent with the mechanisms that we consider here.

In the correlated state, spike train correlations in the recurrent network are essentially inherited from correlations between spike trains in the external population. Hence, the $\mathcal{O}(1)$ correlations realized by this mechanism require the presence of another local network with $\mathcal{O}(1)$ correlations. This raises the question of where the $\mathcal{O}(1)$ correlations are originally generated. One possibility is that they could be generated in a presynaptic cortical area or layer through the alternative mechanisms discussed in the previous paragraph. Another possibility is that they originate from a network that is not in the balanced state at all. Non-balanced networks can easily achieve $\mathcal{O}(1)$ spike train correlations simply from overlapping synaptic projections. While cortical circuits are commonly believed to operate in a balanced state, correlations could originate in thalamus, retina, or other sub-cortical neural populations then eventually propagate to cortex.

The cancellation between variances of covariances observed empirically in Fig. 4F is, 618 to our knowledge, a novel observation, but we were unable to derive it analytically. 619 Path integral approaches have recently been applied to compute variances of covariances 620 in linear network models with uncorrelated external input [59] ($\langle \boldsymbol{X}, \boldsymbol{X} \rangle = 0$), and could 621 potentially be extended to the networks considered here. This previous work derives a 622 simple relationship between a network's criticality and a parameter, Δ , that represents 623 the ratio between the variance of covariances and the mean spike count variance. 624 Specifically, they showed that in their networks, the eigenvalue spectrum of the network 625 dynamics is given to leading order in N by $\lambda_{max} = \sqrt{1 - \sqrt{1/(1 + N\Delta^2)}}$. In our model, 626

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

> where external input is correlated, the variance of covariances and spike count variances both appear to be $\mathcal{O}(1)$, so that $\Delta \sim \mathcal{O}(1)$ and their expression for λ_{max} would tend to zero regardless of the network's dynamical state. We interpret this to imply that their expression is not accurate for networks with correlated external input. Future work should consider the possibility of extending their analysis of criticality to such networks. (33)

Our mean-field analysis applies to networks of integrate-and-fire neuron models, which are arguably more biologically realistic than networks of binary model neurons that are often used for mean-field analysis of neuronal networks. Binary neuron networks are appealing due to the mathematical tractability of their mean-field analysis, but our work demonstrates that integrate-and-fire networks are similarly tractable, calling into question the utility and appeal of binary network models.

Two unproven assumptions underly our mean-field analysis of the correlated state. The first assumption is that neural transfer is $\mathcal{O}(1)$ (Eq. (8) and surrounding discussion). The second assumption is that individual connection strengths are not strongly correlated with individual CSD values so that the step from Eq. (26) to (27) is valid when ignoring smaller order terms. These assumptions are made in other work, even if not stated explicitly. We have been unable to prove these assumptions rigorously for the model studied here, leaving an open problem for future work.

We showed that linear approximations to spike train covariance developed for small, sparsely coupled networks [10–12] can also be accurate for large, densely connected balanced networks (Fig. 5). However, their usefulness is limited by the need to invert large, ill-conditioned matrices and to approximate the susceptibility functions of individual neurons. The simpler equations we derived for mean-field spike train CSDs (Eqs. (7) and (14)) do not have these problems. Moreover, while linear response approximations require that neural transfer of input is approximately linear, our mean-field derivations did not depend on this assumption. Recent work has called for looking beyond linear analysis of neuronal networks [60]. Our analysis shows that, even in networks where neural transfer of inputs is nonlinear, linear mean-field analysis could still be accurate and useful.

In summary, we showed that correlations in balanced networks can be caused by feedforward input from a population of neurons with correlated spike trains, defining the "correlated state" that is quantitatively captured by a linear mean-field theory. In contrast to other mechanisms of correlation in balanced networks, the correlated state predicts a precise balance between the fluctuations in excitatory and inhibitory synaptic input to individual neuron pairs, consistent with some in vivo recordings [22].

Materials and Methods

Details for the derivation of mean-field CSDs.

We now provide details in the derivations of Eqs. (6) and (14), which can both be written as

$$\langle \boldsymbol{S}, \boldsymbol{S} \rangle = \frac{1}{N} W^{-1} \langle \boldsymbol{X}, \boldsymbol{X} \rangle W^{-*} + \frac{1}{N} C_0 + o(\langle \boldsymbol{X}, \boldsymbol{X} \rangle / N)$$
(25)

where $o(\langle \boldsymbol{X}, \boldsymbol{X} \rangle / N)$ scales smaller than $\langle \boldsymbol{X}, \boldsymbol{X} \rangle / N$ as $N \to \infty$ and where $C_0 \sim \mathcal{O}(1)$. Note that in the correlated state, $\langle \boldsymbol{X}, \boldsymbol{X} \rangle \sim \mathcal{O}(N)$ so that the C_0/N term can be absorbed into the $o(\langle \boldsymbol{X}, \boldsymbol{X} \rangle / N)$ term. A sketch of the derivation for the correlated state is given in Results, and the derivation is similar in the asynchronous state. Here, we give the details of this derivation.

We first derive Eq. (5) for $\langle X, X \rangle$ in the asynchronous state, *i.e.* when spike trains in the external population are uncorrelated Poisson processes, so $\langle S_x, S_x \rangle = 0$. In that

663

664

665

666

667

668

669

670

662

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660



case, the CSD between two neurons' external input is given by

$$\begin{split} \langle X_j^a, X_k^b \rangle &= \left\langle \sum_{m=1}^{N_x} J_{jm}^{ax} \eta_x * S_m^x, \sum_{n=1}^{N_x} J_{kn}^{bx} \eta_x * S_n^x \right\rangle \\ &= \sum_{m,n=1}^{N_x} J_{jm}^{ax} J_{kn}^{bx} \widetilde{\eta}_x \widetilde{\eta}_x^* \langle S_m^x, S_n^x \rangle \end{split}$$

for a, b = e, i where we have used the fact that the cross-spectral operator, $\langle \cdot, \cdot \rangle$ is a Hermitian operator. Since external spike trains are uncorrelated Poisson processes, $\langle S_m^x, S_n^x \rangle = 0$ when $m \neq n$ and $\langle S_m^x, S_m^x \rangle = r_x$. Therefore, we can rewrite the equation above as

$$\langle X_j^a, X_k^b \rangle = |\widetilde{\eta}_x|^2 r_x \sum_{m=1}^{N_x} J_{jm}^{ax} J_{km}^{bx}.$$

From Eq. (2), the expectation of the summand in this equation is

$$E[J_{jm}^{ax}J_{km}^{bx}] = \frac{p_{ax}p_{bx}j_{ax}j_{bx}}{N}.$$

Hence, taking expectations across the N_x elements of the sum and the coefficient in front of the sum gives 680

$$X_a, X_b \rangle = q_x j_{ax} j_{bx} p_{ax} p_{bx} |\tilde{\eta}_x|^2 r_x$$

which, in matrix form, is equivalent to Eq. (5).

<

We next derive Eq. (10). Noting that recurrent input to neuron j in population a = e, i is composed of excitatory and inhibitory components, $R_j^a(t) = E_j^a(t) + I_j^a(t)$, we have have

$$\langle R_j^a, U_k^b \rangle = \langle E_j^a, U_k^b \rangle + \langle I_j^a, U_k^b \rangle$$

where we can compute

$$\langle E_j^a, U_k^b \rangle = \left\langle \sum_{m=1}^{N_e} J_{jm}^{ae} \eta_e * S_m^e, U_k^b \right\rangle \tag{26}$$

$$= \widetilde{\eta}_e \sum_{m=1}^{N_e} J_{jm}^{ae} \langle S_m^e, U_k^b \rangle \tag{27}$$

$$= \tilde{\eta}_e \sum_{m \neq k} J_{jm}^{ae} \langle S_m^e, U_k^b \rangle + \tilde{\eta}_e J_{jk}^{ae} \langle S_k^e, U_k^b \rangle.$$
⁽²⁸⁾

Taking expectation over j and k as above gives

$$\langle E_a, U_b \rangle = \sqrt{N} p_{ae} j_{ae} q_e \tilde{\eta}_e \langle S_e, U_b \rangle + \mathcal{O}(\operatorname{avg}_{k,b} \langle S_k^e, U_k^b \rangle / \sqrt{N}).$$

where $\mathcal{O}(\operatorname{avg}_{k,b}\langle S_k^e, U_k^b \rangle / \sqrt{N})$ accounts for the diagonal terms not counted in the definition of $\langle S_e, U_b \rangle$. Note that this step requires us to assume that individual values of the random variable, J_{jm}^{ae} , are not strongly correlated with individual values of $\langle S_m^e, U_k^b \rangle$, so that the expectation of their product can be replaced by the product of their expectations. This assumption is implicit in derivations in other studies [29, 38, 59], even though it is never proven and often not made explicit.

Repeating this calculation for $\langle I_i^a, U_k^b \rangle$ and putting them together gives the average 692

$$\langle R_a, U_b \rangle = \sqrt{N} (w_{ae} \langle S_e, U_b \rangle + w_{ai} \langle S_i, U_b \rangle) + \mathcal{O}(\operatorname{avg}_{k,b,c} \langle S_k^c, U_k^b \rangle / \sqrt{N} \rangle)$$

673

678

681

685

686

687

688

689

690

In matrix form, this becomes

SUBMISSION

$$\langle \boldsymbol{R}, \boldsymbol{U} \rangle = \sqrt{N} W \langle \boldsymbol{S}, \boldsymbol{U} \rangle + \mathcal{O}(\operatorname{avg}_{k \ b \ c} \langle S_k^c, U_k^b \rangle / \sqrt{N} \rangle).$$

An identical calculation, replacing \boldsymbol{R} with \boldsymbol{X} and \boldsymbol{S} with S_x , gives

$$\langle \boldsymbol{X}, \boldsymbol{U} \rangle = \sqrt{N} W_x \langle S_x, \boldsymbol{U} \rangle + \mathcal{O}(\operatorname{avg}_{k,b} \langle S_k^x, U_k^b \rangle / \sqrt{N} \rangle)$$

For the correction terms, $\mathcal{O}(\operatorname{avg}_{k,b,c}\langle S_k^c, U_k^b \rangle / \sqrt{N} \rangle)$ and $\mathcal{O}(\operatorname{avg}_{k,b}\langle S_k^x, U_k^b \rangle / \sqrt{N} \rangle)$, to contribute at largest order in N, it needs to be true that

$$\operatorname{avg}_{k,b}\langle S_k^c, U_k^b \rangle \ge N\mathcal{O}(\langle \boldsymbol{S}, \boldsymbol{U} \rangle)$$

for c = e, i, or x. In the correlated state, this is never the case, so the correction term can be ignored, giving Eqs. (10) and (11). In the asynchronous state, $\operatorname{avg}_{k,b}\langle S_k^b, S_k^b \rangle \sim \mathcal{O}(1)$ since spike train power spectral densities are $\mathcal{O}(1)$ due to intrinsically generated variability, but $\langle \mathbf{S}, \mathbf{S} \rangle \sim \mathcal{O}(1/N)$. This causes the power spectral densities, *i.e.* auto-correlations, to contribute to correlated variability in the network at the largest order in N and ultimately leads to the presence of the C_0 term in Eq. (6) where $(1/N)C_0$ represents mean-field CSDs that would be obtained in the absence of external input correlations, $\langle \mathbf{X}, \mathbf{X} \rangle = 0$ (see other work [38, 39, 44] for an in-depth

Derivation of linear response approximation to pairwise spike train CSDs.

We next give a derivation of Eq. (18) from Eqs. (16) and (17). Similar derivations have previously been given for integrate-and-fire networks [11, 12] and other models [10, 13, 44, 59, 61]. First compute

$$\begin{split} \langle \vec{X}, \vec{T} \rangle &= \langle \vec{X}, \vec{R} \rangle + \langle \vec{X}, \vec{X} \rangle \\ &= \langle \vec{X}, \vec{S} \rangle W^* + \langle \vec{X}, \vec{X} \rangle \\ &\approx \langle \vec{X}, \vec{T} \rangle \widetilde{A}^* W^* + \langle \vec{X}, \vec{X} \rangle \end{split}$$

from Eqs. (16) and (17). This can be solved for

treatment of intrinsically generated correlations).

$$\langle \vec{X}, \vec{T} \rangle = \langle \vec{X}, \vec{X} \rangle (Id - W\widetilde{A})^{-*}.$$

Similarly, compute

$$\begin{split} \langle \vec{T}, \vec{T} \rangle &= \langle \vec{R}, \vec{T} \rangle + \langle \vec{X}, \vec{T} \rangle \\ &= W \langle \vec{S}, \vec{T} \rangle + \langle \vec{X}, \vec{T} \rangle \\ &\approx W \widetilde{A} \langle \vec{T}, \vec{T} \rangle + \langle \vec{X}, \vec{X} \rangle (Id - W \widetilde{A})^{-*} \end{split}$$

which can be solved to obtain

$$\langle \vec{T}, \vec{T} \rangle = (Id - W\widetilde{A})^{-1} \langle \vec{X}, \vec{X} \rangle (Id - W\widetilde{A})^{-*}.$$

Finally, making the substitution $\langle \vec{S}, \vec{S} \rangle = A \langle \vec{T}, \vec{T} \rangle A^*$, which follows from Eqs. (17), ⁷¹⁴ gives Eq. (18). ⁷¹⁵

712

713

711

695

693

694

696

697

698

699

700

701

702

703

704

705

706

Generation of correlated spike trains for external inputs.

To generate correlated, Poisson spike trains for the external population in the correlated 717 state we used the multiple interaction process (MIP) method [62] with jittering. 718 Specifically, we generated one shared "mother" process with firing rate $r_m = r_x/c$. 719 Then, for each of the N_x "daughter" processes, we randomly kept each spike in the 720 mother process with probability c. As a result, each daughter process is a Poisson 721 process with firing rate $cr_m = r_x$ and a proportion of c of the spikes are shared between 722 any two daughter processes. To get rid of perfect synchrony between the daughter 723 processes, we jittered each spike time in each daughter process by a normally 724 distributed random variable with mean zero and standard deviation $\tau_c = 5$ ms. Upon 725 jittering, the daughter processes remain Poisson and the resulting CSD between 726 daughter processes is given by Eq. (15). Spike count correlations between the daughter 727 processes over large time windows are exactly c. The daughter processes were used as 728 the spike trains, $S_{i}^{*}(t)$ in the external population in the correlated state. See [62] for a 729 deeper analysis of this algorithm. 730

Parameters for simulations

All connection probabilities were $p_{ab} = 0.1$ for a = e, i and b = e, i, x. Synaptic timescales were $\tau_e = 8 \text{ms}$, $\tau_i = 4 \text{ms}$, and $\tau_x = 10 \text{ms}$. The firing rate of the external population was $r_x = 10\text{Hz}$ and, in the correlated state, the correlation was c = 0.1 with a jitter of $\tau_c = 5 \text{ms}$. All covariances and correlations were computed by counting spikes or integrating continuous processes over a window of length 250 ms. Membrane capacitance, C_m , is arbitrary so we report all current-based parameters in relation to C_m . For convenience, one can therefore set $C_m = 1$. Unscaled connection strengths were $j_{ee}/C_m = 25 \text{mV}$, $j_{ei}/C_m = -150 \text{mV}$, $j_{ie}/C_m = 112.5 \text{mV}$, $j_{ii}/C_m = -250 \text{mV}$, $j_{ex}/C_m = 180 \text{mV}$, and $j_{ix}C_m = 135 \text{mV}$. Note that j_{ab} was scaled by \sqrt{N} to produce the true connection strengths, as indicated in Results. Neuron parameters are $g_L = C_m/15$, $E_L = -72 \text{mV}$, $V_{th} = -50 \text{mV}$, $V_{re} = -75 \text{mV}$, $V_{lb} = -100 \text{mV}$, $\Delta_T = 1 \text{mV}$, and $V_T = -55 \text{mV}$. Synaptic currents in figures are reported in units $C_m V/s$. Covariances between synaptic currents are computed between integrals of the currents (see Eq. (3) and surrounding discussion), so the covariances have units $C_m^2 MV^2$.

Details of computer simulations.

All simulations and numerical computations were performed on a MacBook Pro running 747 OS X 10.9.5 with a 2.3 GHz Intel Core i7 processor. All simulations were written in 748 Matlab (Matlab R 2018a, MathWorks). The differential equations defining the neuron 749 model were solved using a forward Euler method with time step 0.1ms. Statistics in 750 Figs. 1D, 2C,D, and 4C,D were computed from a simulation of duration 50s. Statistics 751 in Figs. 1E-I, 2E–I, and 4A, B, E, F were computed by repeating a simulation of duration 752 50s over ten trials for each value of N, then averaging over trials. For each trial, 753 network connectivity was generated with a different random seed, so the statistics are 754 averaged over time and over realizations of the "quenched" variability of network 755 connectivity. Statistics in Fig. 5 were computed by repeating a simulation of duration 756 100s for 50 trials, with network connectivity the same for each trial. Statistics were then 757 averaged over trials. Gains were estimated by fitting a rectified quadratic function to 758 the relationship between all neurons' firing rates and mean total inputs $(r_i \text{ and } \overline{T}_i)$, 759 then computing the derivative of the fitted quadratic at the input value for each neuron. 760 The same approach was used in previous work [38, 51] to estimate a mean-field gain. 761 Matlab files to produce all figures are available from the XXX database (accession 762 number(s) XXX, XXX). 763

716

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

Acknowledgments

We thank Michael S. Okun for helpful comments on a draft of the manuscript and for his contribution to collecting the data used in Fig. 3A. RR was supported by National Science Foundation (https://www.nsf.gov/) grant numbers DMS-1517828, DMS-1654268, and DBI-1707400. IL was supported by Deutsche Forschungsgemeinschaft Sonderforschungsbereiche (www.dfg.de/sfb) grant number 1089, Israel Science Foundation (https://www.isf.org.il/) grant numbers 326/07 and 1539/17, and Minerva (http://www.minerva.mpg.de/weizmann/). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

References

1.	Moreno-Bote R, Renart A, Parga N. Theory of input spike auto-and	775
	cross-correlations and their effect on the response of spiking neurons. Neural	776
	computation. $2008;20(7):1651-1705.$	777
2.	Pouget A Beck JM Ma WJ Latham PE Probabilistic brains: knowns and	778

- 2. Pouget A, Beck JM, Ma WJ, Latham PE. Probabilistic brains: knowns and unknowns. Nature neuroscience. 2013;16(9):1170–1178.
- 3. Shamir M. Emerging principles of population coding: in search for the neural code. Current opinion in neurobiology. 2014;25:140–148.
- Tetzlaff T, Rotter S, Stark E, Abeles M, Aertsen A, Diesmann M. Dependence of neuronal correlations on filter characteristics and marginal spike train statistics. Neural Comput. 2008;20(9):2133–84. doi:101162/neco200805-07-525.
- Tetzlaff T, Helias M, Einevoll GT, Diesmann M. Decorrelation of neural-network activity by inhibitory feedback. PLoS Comput Biol. 2012;8(8):e1002596.
- Doiron B, Litwin-Kumar A, Rosenbaum R, Ocker GK, Josić K. The mechanics of state-dependent neural correlations. Nature Neuroscience. 2016;19(3):383–393.
- Ocker GK, Hu Y, Buice MA, Doiron B, Josić K, Rosenbaum R, et al. From the statistics of connectivity to the statistics of spike times in neuronal networks. Current opinion in neurobiology. 2017;46:109–119.
- Lindner B, Doiron B, Longtin A. Theory of oscillatory firing induced by spatially correlated noise and delayed inhibitory feedback. Phys Rev E. 2005;72(6):061919.
- Ostojic S, Brunel N, Hakim V. How connectivity, background activity, and synaptic properties shape the cross-correlation between spike trains. Journal of Neuroscience. 2009;29(33):10234–10253.
- Pernice V, Staude B, Cardanobile S, Rotter S. How structure determines correlations in neuronal networks. PLoS Comput Biol. 2011;7(5):e1002059.
 doi:10.1371/journal.pcbi.1002059.
- 11. Trousdale J, Hu Y, Shea-Brown E, Josić K. Impact of network structure and cellular response on spike time correlations. PLoS Comput Biol.
 800

 2012;8(3):e1002408.
 802
- 12. Hu Y, Trousdale J, Josić K, Shea-Brown E. Motif statistics and spike correlations in neuronal networks. J Stat Mech. 2013;2013(03):P03012.

764

765

766

767

768

769

770

771

772

773

774

779

780



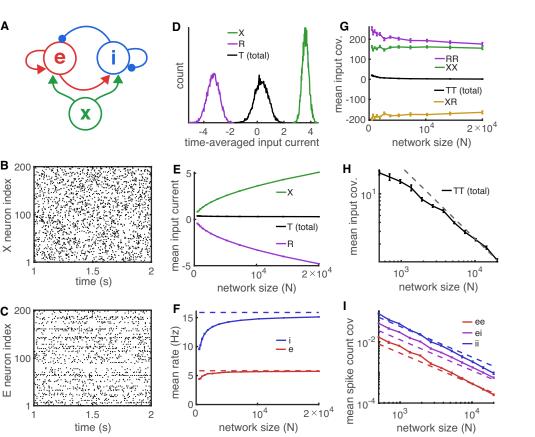
13.	Grytskyy D, Tetzlaff T, Diesmann M, Helias M. A unified view on weakly correlated recurrent networks. Front Comput Neurosci. 2013;7(October):131. doi:10.3389/fncom.2013.00131.	805 806 807
14.	Dummer B, Wieland S, Lindner B. Self-consistent determination of the spike-train power spectrum in a neural network with sparse connectivity. Frontiers in computational neuroscience. 2014;8:104.	808 809 810
15.	Ko H, Hofer SB, Pichler B, Buchanan KA, Sjöström PJ, Mrsic-Flogel TD. Functional specificity of local synaptic connections in neocortical networks. Nature. 2011;473(7345):87–91. doi:101038/nature09880.	811 812 813
16.	Fino E, Yuste R. Dense inhibitory connectivity in neocortex. Neuron. 2011;69(6):1188–1203.	814 815
17.	Levy RB, Reyes AD. Spatial profile of excitatory and inhibitory synaptic connectivity in mouse primary auditory cortex. J Neurosci. 2012;32(16):5609–5619.	816 817 818
18.	Oswald AM, Doiron B, Rinzel J, Reyes AD. Spatial profile and differential recruitment of GABAB modulate oscillatory activity in auditory cortex. J Neurosci. 2009;29(33):10321–10334.	819 820 821
19.	Shu Y, Hasenstaub A, McCormick DA. Turning on and off recurrent balanced cortical activity. Nature. 2003;423(6937):288–293.	822 823
20.	Wehr M, Zador AM. Balanced inhibition underlies tuning and sharpens spike timing in auditory cortex. Nature. 2003;426(6965):442–446.	824 825
21.	Haider B, Duque A, Hasenstaub AR, McCormick DA. Neocortical network activity in vivo is generated through a dynamic balance of excitation and inhibition. J Neurosci. 2006;26(17):4535–4545.	826 827 828
22.	Okun M, Lampl I. Instantaneous correlation of excitation and inhibition during ongoing and sensory-evoked activities. Nat Neurosci. 2008;11(5):535–537.	829 830
23.	Dorrn AL, Yuan K, Barker AJ, Schreiner CE, Froemke RC. Developmental sensory experience balances cortical excitation and inhibition. Nature. 2010;465(7300):932–936.	831 832 833
24.	Sun YJ, Wu GK, Liu Bh, Li P, Zhou M, Xiao Z, et al. Fine-tuning of pre-balanced excitation and inhibition during auditory cortical development. Nature. 2010;465(7300):927–931.	834 835 836
25.	Zhou M, Liang F, Xiong XR, Li L, Li H, Xiao Z, et al. Scaling down of balanced excitation and inhibition by active behavioral states in auditory cortex. Nat Neurosci. 2014;17(6):841–50. doi:101038/nn3701.	837 838 839
26.	Petersen PC, Vestergaard M, Jensen KHR, Berg RW. Premotor spinal network with balanced excitation and inhibition during motor patterns has high resilience to structural division. J Neurosci. 2014;34(8):2774–84. doi:101523/JNEUROSCI3349-132014.	840 841 842 843
27.	van Vreeswijk C, Sompolinsky H. Chaotic balanced state in a model of cortical circuits. Neural Comput. 1998;10(6):1321–1371.	844 845
28.	van Vreeswijk C, Sompolinsky H. Chaos in neuronal networks with balanced excitatory and inhibitory activity. Science. 1996;274(5293):1724–1726.	846 847

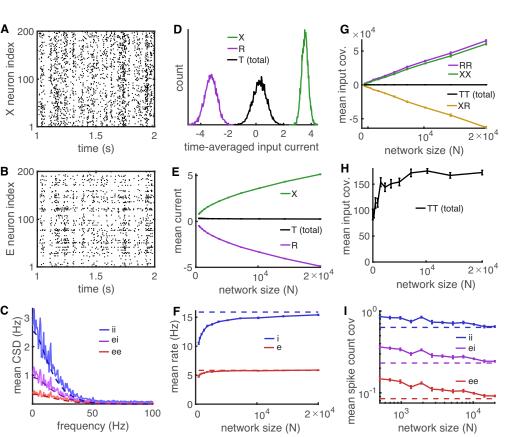


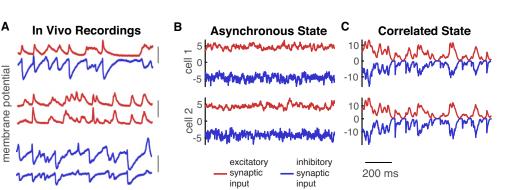
29.	Renart A, de La Rocha J, Bartho P, Hollender L, Parga N, Reyes A, et al. The Asynchronous State in Cortical Circuits. Science. 2010;327(5965):587–590.	848 849
30.	Ecker A, Berens P, Keliris G, Bethge M, Logothetis N, Tolias A. Decorrelated Neuronal Firing in Cortical Microcircuits. Science. 2010;327(5965):584–587.	850 851
31.	Cohen MR, Kohn A. Measuring and interpreting neuronal correlations. Nat Neurosci. 2011;14(7):811–819.	852 853
32.	Smith MA, Jia X, Zandvakili A, Kohn A. Laminar dependence of neuronal correlations in visual cortex. J Neurophysiol. 2013;109(4):940–947.	854 855
33.	Ecker AS, Berens P, Cotton RJ, Subramaniyan M, Denfield GH, Cadwell CR, et al. State dependence of noise correlations in macaque primary visual cortex. Neuron. 2014;82(1):235–248.	856 857 858
34.	Tan AY, Chen Y, Scholl B, Seidemann E, Priebe NJ. Sensory stimulation shifts visual cortex from synchronous to asynchronous states. Nature. 2014;509(7499):226.	859 860 861
35.	McGinley MJ, Vinck M, Reimer J, Batista-Brito R, Zagha E, Cadwell CR, et al. Waking state: rapid variations modulate neural and behavioral responses. Neuron. 2015;87(6):1143–1161.	862 863 864
36.	Mochol G, Hermoso-Mendizabal A, Sakata S, Harris KD, de la Rocha J. Stochastic transitions into silence cause noise correlations in cortical circuits. Proc Natl Acad Sci USA. 2015;112(11):201410509. doi:101073/pnas1410509112.	865 866 867
37.	Wimmer K, Compte A, Roxin A, Peixoto D, Renart A, de la Rocha J. The dynamics of sensory integration in a hierarchical network explains choice probabilities in MT. Nat Commun. 2015;6:1–13. doi:101038/ncomms7177.	868 869 870
38.	Rosenbaum R, Smith MA, Kohn A, Rubin JE, Doiron B. The spatial structure of correlated neuronal variability. Nature Neurosci. 2017;20(1):107.	871 872
39.	Darshan R, van Vreeswijk C, Hansel D. How strong are correlations in strongly recurrent neuronal networks? 2018;(1):1–22. doi:10.1101/274480.	873 874
40.	Rosenbaum R, Doiron B. Balanced networks of spiking neurons with spatially dependent recurrent connections. Phys Rev X. 2014;4(2):021039. doi:101103/PhysRevX4021039.	875 876 877
41.	Landau ID, Egger R, Dercksen VJ, Oberlaender M, Sompolinsky H. The impact of structural heterogeneity on excitation-inhibition balance in cortical networks. Neuron. 2016;92(5):1106–1121.	878 879 880
42.	Pyle R, Rosenbaum R. Highly connected neurons spike less frequently in balanced networks. Phys Rev E. 2016;93(4):040302.	881 882
43.	Monteforte M, Wolf F. Dynamic flux tubes form reservoirs of stability in neuronal circuits. Phys Rev X. 2012;2(4):041007.	883 884
44.	Helias M, Tetzlaff T, Diesmann M. The correlation structure of local neuronal networks intrinsically results from recurrent dynamics. PLoS Comput Biol. 2014;10(1):e1003428.	885 886 887
45.	Rosenbaum RJ, Trousdale J, Josić K. Pooling and correlated neural activity. Front Comp Neurosci. 2010;4:1–14.	888 889

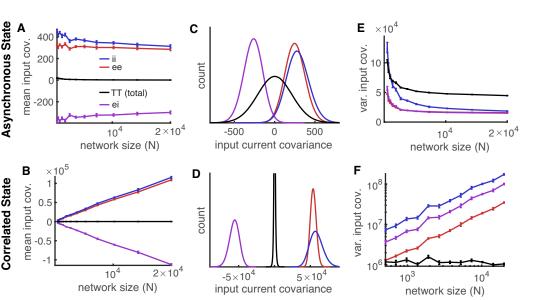


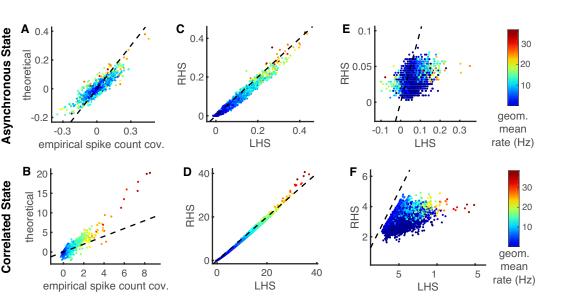
46.	Rosenbaum R, Trousdale J, Josić K. The effects of pooling on spike train correlations. Front Neurosci. 2011;5(58):1–10.	890 891
47.	Hennequin G, Agnes EJ, Vogels TP. Inhibitory Plasticity: Balance, Control, and Codependence. Annu Rev Neurosci. 2017;40(1):557–579. doi:10.1146/annurev-neuro-072116-031005.	892 893 894
48.	Risken H. Fokker-planck equation. In: The Fokker-Planck Equation. Springer; 1996. p. 63–95.	895 896
49.	Lindner B. Effects of noise in excitable systems. Phys Rep. 2004;392(6):321–424. doi:10.1016/j.physrep.2003.10.015.	897 898
50.	Rosenbaum R. A Diffusion Approximation and Numerical Methods for Adaptive Neuron Models with Stochastic Inputs. Front Comput Neurosci. 2016;10(April):1–20. doi:10.3389/fncom.2016.00039.	899 900 901
51.	Ebsch C, Rosenbaum R. Imbalanced amplification: A mechanism of amplification and suppression from local imbalance of excitation and inhibition in cortical circuits. PLoS computational biology. 2018;14(3):e1006048.	902 903 904
52.	Roxin A, Brunel N, Hansel D. Role of delays in shaping spatiotemporal dynamics of neuronal activity in large networks. Physical review letters. 2005;94(23):238103.	905 906
53.	Kriener B, Helias M, Rotter S, Diesmann M, Einevoll GT. How pattern formation in ring networks of excitatory and inhibitory spiking neurons depends on the input current regime. Frontiers in Computational Neuroscience. 2014;7(187). doi:10.3389/fncom.2013.00187.	907 908 909 910
54.	Ostojic S. Two types of asynchronous activity in networks of excitatory and inhibitory spiking neurons. Nature neuroscience. 2014;17(4):594.	911 912
55.	Keane A, Gong P. Propagating waves can explain irregular neural dynamics. Journal of Neuroscience. 2015;35(4):1591–1605.	913 914
56.	Pyle R, Rosenbaum R. Spatiotemporal dynamics and reliable computations in recurrent spiking neural networks. Physical Rev Lett. 2017;118(1):018103.	915 916
57.	Huang C, Ruff D, Pyle R, Rosenbaum R, Cohen M, Doiron BD. Circuit-based models of shared variability in cortical networks. bioRxiv. 2017; p. 217976.	917 918
58.	Malina KCK, Mohar B, Rappaport AN, Lampl I. Local and thalamic origins of correlated ongoing and sensory-evoked cortical activities. Nat Comm. 2016;7:12740.	919 920 921
59.	Dahmen D, Grün S, Diesmann M, Helias M. Two types of criticality in the brain. 2017;(1):1–15.	922 923
60.	Herfurth T, Tchumatchenko T. How linear response shaped models of neural circuits and the quest for alternatives. Current opinion in neurobiology. 2017;46:234–240.	924 925 926
61.	Gardiner CW. Handbook of stochastic methods. vol. 3. Springer, Berlin; 1985.	927
62.	Kuhn A, Aertsen A, Rotter S. Higher-order statistics of input ensembles and the response of simple model neurons. Neural Comput. 2003;15(1):67–101. doi:10.1162/089976603321043702.	928 929 930

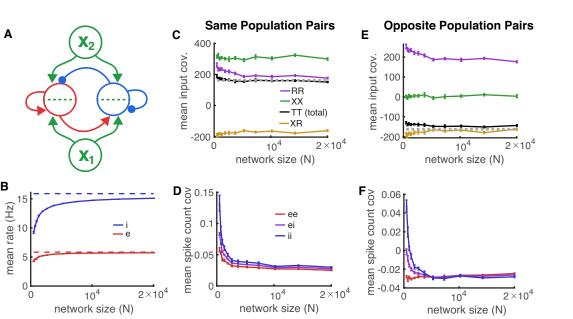




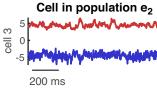








Cells in population e₁



excitatory synaptic input input input inhibitory synaptic input