

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23

Multistability in neural systems with random cross-connections

Jordan Breffle¹, Subhadra Mokashe¹, Siwei Qiu^{2,4}, Paul Miller^{1,2,3*}

¹Neuroscience Program, ²Volen National Center for Complex Systems, and

³Department of Biology, Brandeis University, 415 South St, Waltham, MA 02454

⁴Current address: Department of Neurology, Cedars-Sinai Medical Center, Los Angeles, CA, USA

*Corresponding Author: pmiller@brandeis.edu

ORCID ids: Jordan Breffle, [0000-0001-5793-4427](https://orcid.org/0000-0001-5793-4427)

Subhadra Mokashe, [0000-0002-5425-2903](https://orcid.org/0000-0002-5425-2903)

Siwei Qiu, [0000-0003-1826-7953](https://orcid.org/0000-0003-1826-7953)

Paul Miller, [0000-0002-9280-000X](https://orcid.org/0000-0002-9280-000X)

Abstract

Neural circuits with multiple discrete attractor states could support a variety of cognitive tasks according to both empirical data and model simulations. We assess the conditions for such multistability in neural systems, using a firing-rate model framework, in which clusters of neurons with net self-excitation are represented as units, which interact with each other through random connections. We focus on conditions in which individual units lack sufficient self-excitation to become bistable on their own. Rather, multistability can arise via recurrent input from other units as a network effect for subsets of units, whose net input to each other when active is sufficiently positive to maintain such activity. In terms of the strength of within-unit self-excitation and standard-deviation of random cross-connections, the region of multistability depends on the firing-rate curve of units. Indeed, bistability can arise with

24 zero self-excitation, purely through zero-mean random cross-connections, if the firing-rate curve rises
25 supralinearly at low inputs from a value near zero at zero input. We simulate and analyze finite systems,
26 showing that the probability of multistability can peak at intermediate system size, and connect with
27 other literature analyzing similar systems in the infinite-size limit. We find regions of multistability with a
28 bimodal distribution for the number of active units in a stable state. Finally, we find evidence for a log-
29 normal distribution of sizes of attractor basins, which can appear as Zipf's Law when sampled as the
30 proportion of trials within which random initial conditions lead to a particular stable state of the system.

31

32 **Keywords:** Attractor basin; mean field; fixed points; bistable.

33

34 **Statements and Declarations.**

35 The work was supported by a grant from the National Institutes of Health, R01 NS104818, by the Swartz
36 Foundation, and by the Neuroscience Graduate Program of Brandeis University.

37

38 **Author Contributions.**

39 All authors contributed to the writing and editing of the manuscript and approved the final version.

40 Simulations were carried out by Jordan Breffle, analysis by Subhadra Mokhashe, Siwei Qiu, and Paul

41 Miller. The first draft was written by Paul Miller, who also conceived of the project.

42

43

44 1. Introduction

45 An extensive literature in neuroscience suggests that neural activity can proceed through sequences of
46 distinct states during sensory processing, motor output, or memory-based decision making (Abeles et
47 al., 1995; Benozzo et al., 2021; Escola et al., 2011; Jones et al., 2007; La Camera et al., 2019; Mazzucato
48 et al., 2015; Miller, 2016; Morcos & Harvey, 2016; Rainer & Miller, 2000; Seidemann et al., 1996). The
49 distinct states are revealed as patterns of neural activity that remain relatively stable for durations much
50 longer than those of the rapid transitions between states. Models of the underlying circuitry assume the
51 states correspond to fixed points (or the remnants of fixed points) of the system (Ballintyn et al., 2019;
52 La Camera et al., 2019; Mazzucato et al., 2019; Miller, 2013; Miller & Katz, 2010; Rabinovich et al., 2001;
53 Rabinovich et al., 2014; Recanatesis et al., 2022; Taylor et al., 2022) with the itinerancy from fixed point
54 to fixed point known as latching dynamics (Boboeva et al., 2021; Lerner et al., 2012, 2014; Lerner &
55 Shriki, 2014; Linkerhand & Gros, 2013; Russo & Treves, 2012; Song et al., 2014; Treves, 2005).
56 Transitions between fixed points can be due to their inherent instability when they are saddle points.
57 Otherwise, in networks where a reduced model of the system possesses multiple stable fixed points,
58 transitions arise from one or more of (1) an external stimulus, (2) noise fluctuations, or (3) the drift of a
59 slow variable which impacts a parameter in the reduced model causing it to cross a bifurcation point.
60 Since the number of stable fixed points becomes a key indicator of the potential information processing
61 or memory capacity of the network, it is important to understand the conditions under which a system
62 possess multiple stable fixed points.

63 Here we use firing-rate models (Wilson & Cowan, 1973), in which each unit represents a cluster
64 or assembly of similarly responsive neurons with stronger connections within each cluster as observed
65 in some cortical circuits (Perin et al., 2011; Song et al., 2005). Such assemblies can arise in response to a
66 lifetime of stimuli via Hebbian plasticity (Hebb, 1949), which increases connection strengths between
67 excitatory neurons with correlated activity (Bourjaily & Miller, 2011; Brunel, 2003). We assume random

68 interactions between such clusters (Stern et al., 2014), representing the result of a history of
69 uncorrelated stimuli.

70 Each isolated stable fixed point in a system is an attractor state, with a basin of attraction
71 determined by the set of initial conditions that result in neural activity settling at (after being “attracted
72 to”) the fixed point. Systems with many such attractor states have provided the framework for
73 understanding pattern completion and separation of new inputs following memory encoding of stimuli,
74 since the highly influential work of Hopfield and others (Anishchenko & Treves, 2006; Battaglia & Treves,
75 1998; Hopfield, 1982; Hopfield, 1984; Treves, 1990; Zurada et al., 1996). Indeed, there is abundant
76 evidence of such attractor states in neural circuits (Daelli & Treves, 2010; Fuster, 1973; Goldberg et al.,
77 2004; Golos et al., 2015; Wills et al., 2005), perhaps most obvious to us when an ambiguous stimulus can
78 cause perceptual alternation due to activity flipping between two (quasi-stable) attractor states
79 (Moreno-Bote et al., 2007). However, while the number of stable states in systems such as the Hopfield
80 network (Hopfield, 1982; Hopfield, 1984) have been characterized (Amit et al., 1985a, 1985b; Folli et al.,
81 2016), the connections between units in such networks are correlated (in fact, the connectivity matrix is
82 symmetric), so it is unclear to what extent multiple attractor states would arise in a nonsymmetric
83 random network.

84 Work by others (Stern et al., 2014) showed that when each unit has sufficient self-excitation to
85 become bistable (and therefore become in essence a memory element in of itself) multiple attractor
86 states are possible in a network with non-symmetrically randomly connected units. Such a result is trivial
87 in the limit of zero cross-connection strength, in which case a system of N bistable units possesses 2^N
88 stable states. In the randomly connected system, studied in the large- N limit, increased strength of
89 random cross-connections decreases the number of multistable states, eventually rendering the system
90 chaotic as all fixed points become unstable. With weaker self-connections, the network would be either
91 quiescent or, given sufficient cross-connection strength, chaotic (Sompolinsky et al., 1988).

92 Here we find that such results depend on the form of the input-output function (the firing-rate,
93 or f-I curve) of a neuron. Indeed, if we assume neurons have low firing rates in the absence of input,
94 random non-symmetric cross-connections can lead to multistability, even when individual units have
95 zero self-excitation.

96 In the following sections, we first present simulations showing the types of activity possible and
97 their observed coexistence in networks of up to 1000 randomly coupled units. We then show the phase
98 diagrams in the large- N limit of such systems. Finally, we present results for systems with binary
99 activation functions, for which we develop an alternative mean-field analytic approach that we use for
100 finite- as well as infinite- N systems. Also, given the more rapid simulations when activations are binary,
101 we provide a more thorough analysis of the attractor states in such systems.

102

103 2. Simulations of Finite Networks

104 We simulated networks of N randomly connected firing rate units with response function $f(x)$
105 representing their output to total input, x . The total input, x_i , to the i -th unit is described by the
106 dynamical equation with time constant, τ :

$$\tau \dot{x}_i = -x_i + s f(x_i) + \frac{g}{\sqrt{N}} \sum_{i \neq j} J_{ij} f(x_j) + I_{Global} \quad (1)$$

107 where s and g are parameters that scale the self-connection and cross-connection strengths,
108 respectively, and J_{ij} is the matrix of normalized cross-connection strengths drawn from a normal
109 distribution with zero mean and unit variance. I_{Global} is a constant input that inhibits or excites the
110 whole network (equivalent to a shift in threshold, x_{th}) and is kept at zero unless stated otherwise.

111 We simulated models with distinct single-unit response functions, $f(x)$, in order to assess its
112 role in network dynamics:

113 1) Hyperbolic tangent: $f(x) = \tanh(x)$ and our model is identical to that of (Stern et al., 2014).

114 We also compare more general forms, $f(x) = \tanh\left(\frac{x-x_{th}}{\Delta}\right)$ to connect results to those of the logistic

115 function with less symmetry in the firing rates (*i.e.*, units require net excitatory input to reach half their

116 maximum rate if $x_{th} > 0$).

117 2) Logistic function: $f(x) = \frac{1}{1+e^{-\frac{x-x_{th}}{\Delta}}}$ where x_{th} is a threshold input required for the firing rate of

118 a unit to reach half of its maximum value and Δ is inversely proportional to the steepness of the

119 response function.

120 3) Binary output via the Heaviside function: $f(x) = Heaviside(x - x_{th})$, which is equivalent to

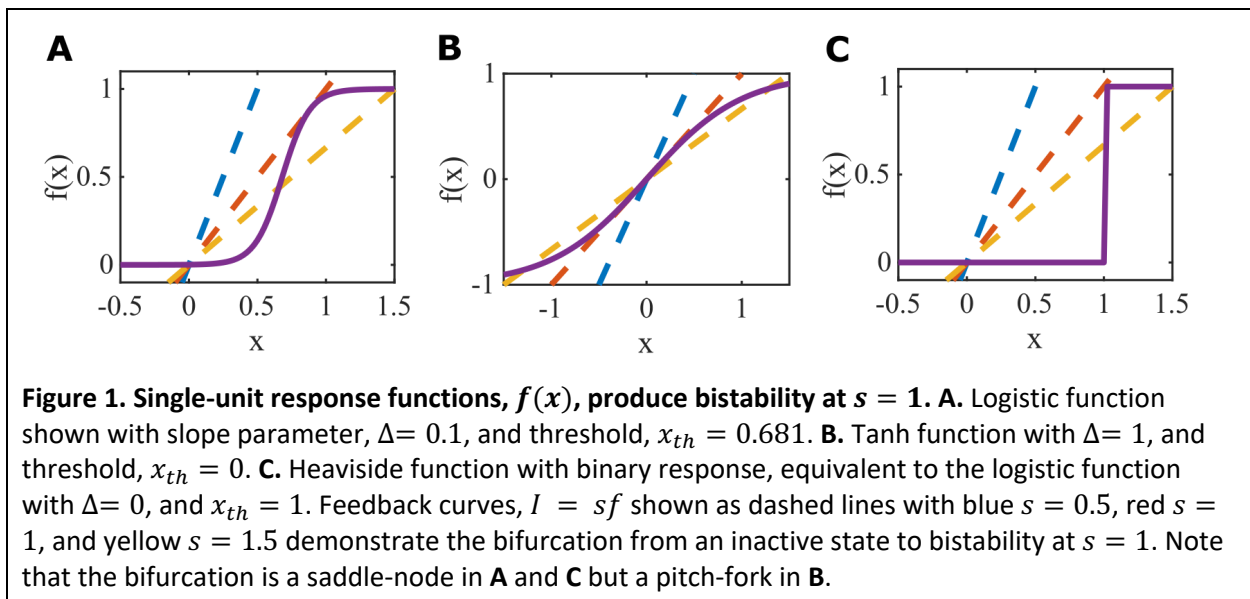
121 the logistic function in the limit $\Delta \rightarrow 0$.

122 In all systems we adjusted the threshold parameter, x_{th} , for a given steepness of response

123 function (*i.e.*, a given value of Δ), such that in the absence of cross-connections ($g = 0$) the system

124 becomes multistable, because each unit is bistable, at $s = 1$. This allows us to compare results across

125 systems with different single-unit response functions, $f(x)$ (Figure 1, see Appendix A).



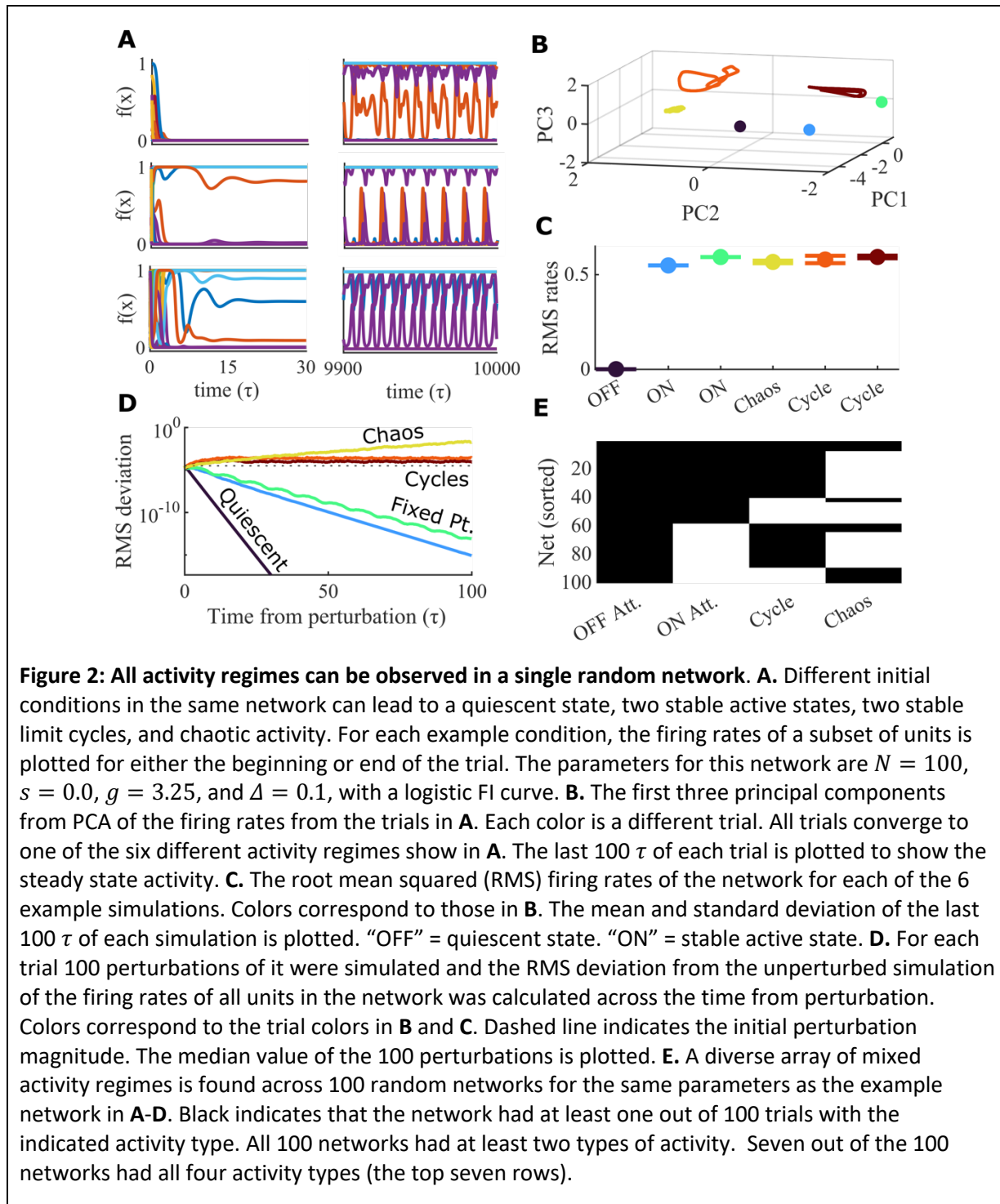
126

127 2.1. Observed forms of simulated network dynamics with logistic function responses

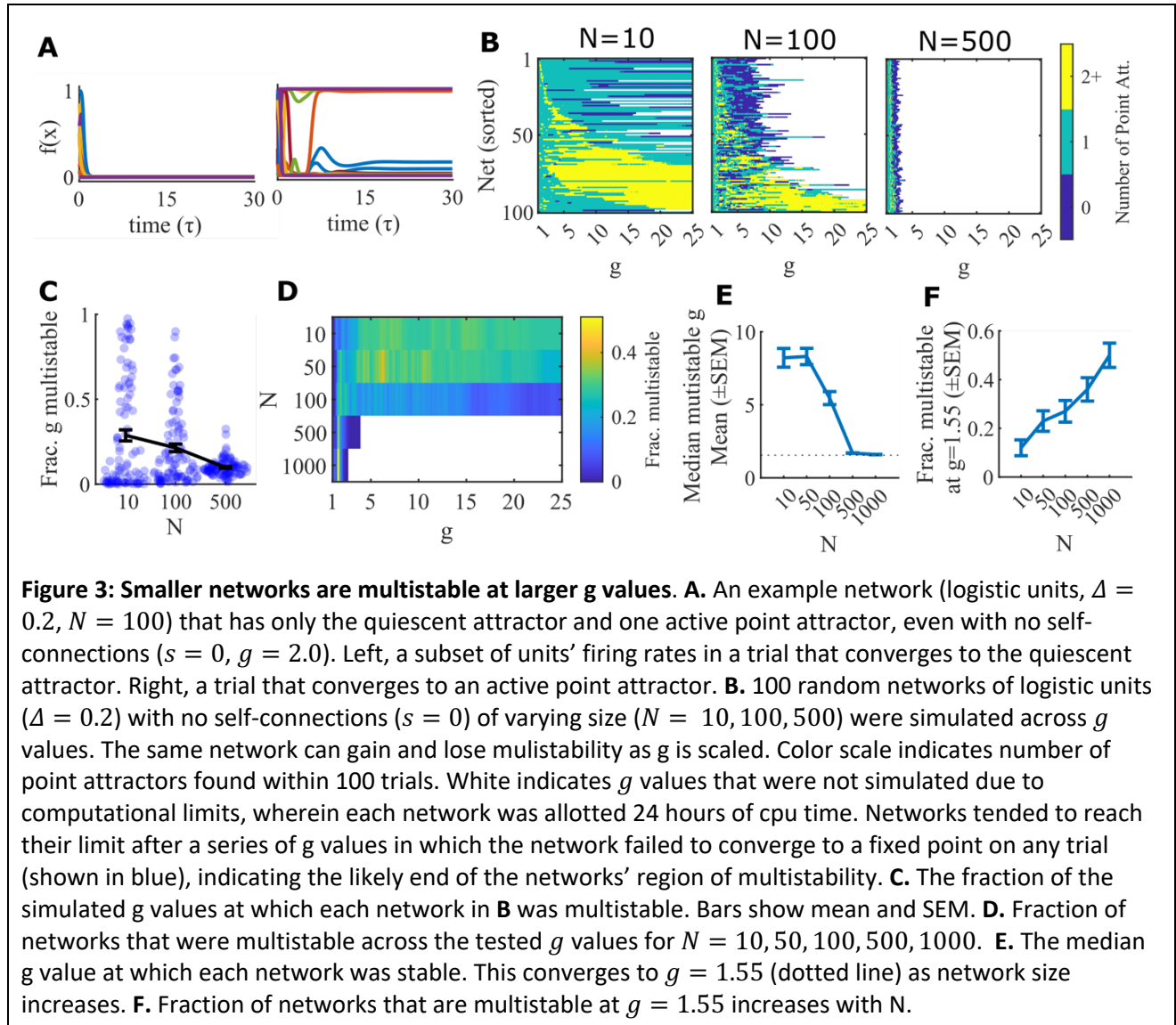
128 We define the network state by its long-term activity, which can be either constant, oscillating, or
129 chaotic. We separate out constant (stable) states into two types: those that include active units and
130 those with only inactive units (and are quiescent). We therefore obtain four labels for final states:
131 quiescence, stable activity, limit cycle, and chaos.

132 In many networks we find, by varying initial conditions, the existence of more than one type of
133 state in a single network. Some networks have multiple forms of all four activity types, such as the
134 example network in Figure 2A-D. This example network of logistic units ($N = 100$, $s = 0$, $g = 3.25$, and
135 $\Delta = 0.1$) has a stable quiescent state, two stable active states, two unique stable limit cycles, and a
136 chaotic attractor. Example trials leading to each of these distinct states in the same network are shown
137 as a subset of units' firing rates (Figure 2A) and in principle component space (Figure 2B). These states
138 have similar root mean squared (RMS) firing rates, except for the quiescent state (Figure 2C).
139 Perturbation analysis confirms the classification of each of the trials in Figure 2A-C (Figure 2D). For each
140 trial, we simulated 100 perturbations and calculated the median RMS deviation of the perturbed
141 simulation from the original simulation. Analysis of the cross connections of the example network in
142 Figure 2A-D shows that it is not an outlier from an expected random network (Supplemental Figure 1),
143 suggesting this combination of mixed activity states may be a common occurrence. We performed the
144 same perturbation-based classification of activity states for 100 different random networks at the same
145 parameter values and found that all networks show at least two forms of activity (Figure 2E).

146 Of particular interest are systems without self-connections, ($s = 0$, such as that shown in Figure
147 2), for which there is a well-established single transition from quiescence to chaos at $g = 1$, when
148 $f(x) = \tanh(x)$ (Sompolinsky et al., 1988). When, instead, we use the logistic function, $f(x) =$
149 $\frac{1}{1+e^{-\frac{x_{th}-x}{\Delta}}}$, which is simply a scaled and shifted transformation of the tanh function to non-negative values
150 of firing rate, we find a richer set of states in our simulations. Perhaps surprising, circuits without self-



151 connections can exhibit multiple stable states: sometimes having only a low activity state (a quiescent
 152 state) with a state of higher net activity (an active state) such as the example network in Figure 3A.



153 Indeed, we find that for a given random instantiation of the connectivity matrix, as we scale all
 154 connections by g , there is, for all 500 total networks tested in Figure 3, some range of connection
 155 strengths for which the network is multistable. Supplemental Figure 2 shows how dynamics beyond the
 156 existence of multistability changes as g is scaled.

157 We wondered whether such states were the results of a finite size effect, so varied the size of
 158 the network (changing N). We found that the range of g over which we see such multistability at $s = 0$
 159 narrows with increased N (Figure 3C), and converges to the same set of values centered on $g = 1.55$

160 (Figure 3D-F). These and other results prompted us to investigate the phase space of the corresponding
161 infinite- N systems via mean field theory and stability analysis (Section 3).

162

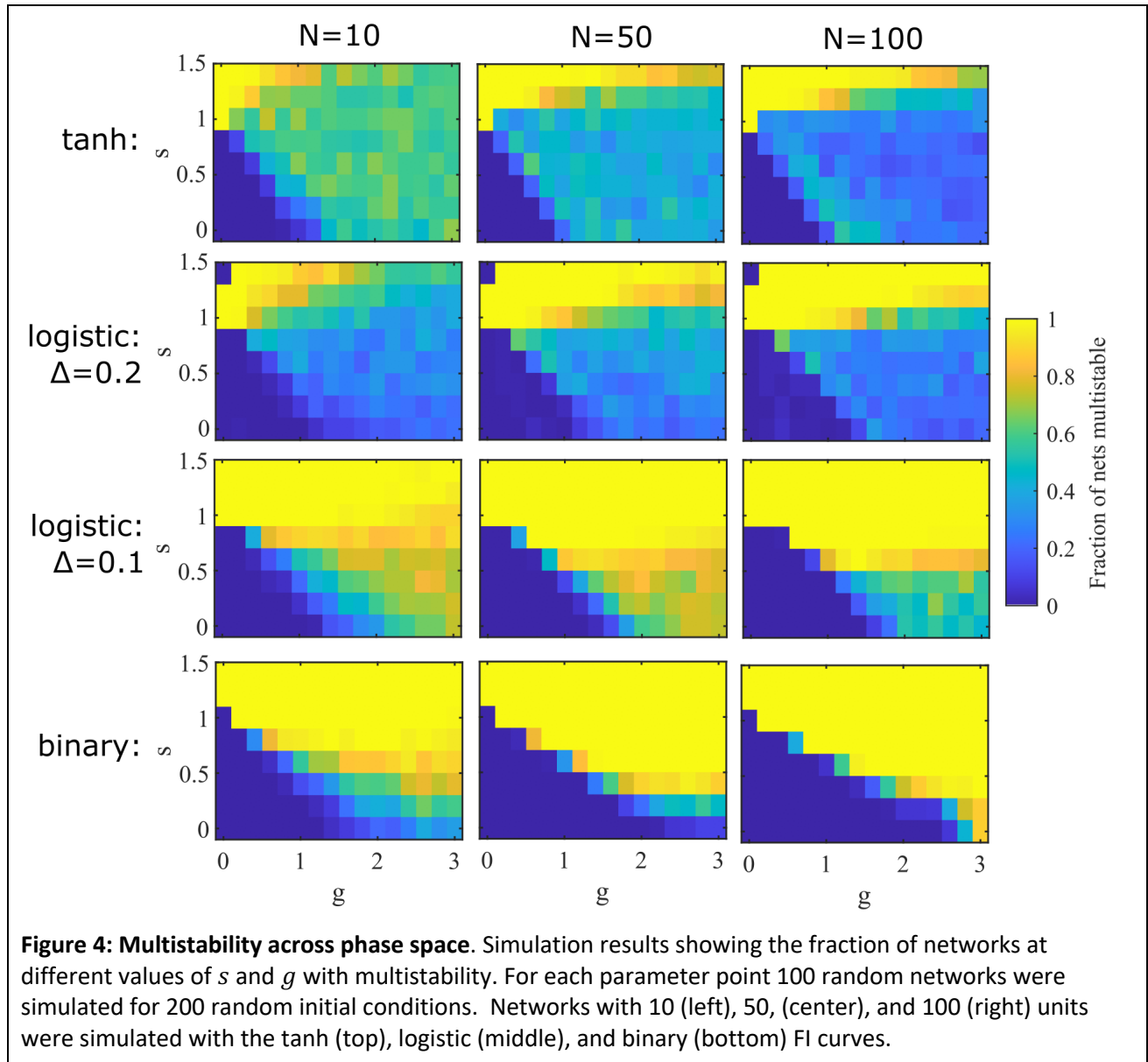
163 2.2. Phase diagram from simulations of finite networks

164 To assess the likelihood of systems reaching a given type of state across phase space, we simulated 100
165 networks for each given set of parameters and commenced simulations of each network from 200 initial
166 conditions. Full details of simulation methods are provided in Appendix 1.

167 In Figure 4 we show that systems with $f(x) = \frac{1}{1+e^{\frac{x_{th}-x}{\Delta}}}$ can have multistability at $s < 1$, such
168 that increasing the cross-coupling strength, g , increases the fraction of multistable networks for large
169 ranges of s and g . The observed multistability at low s that arises with increasing g is most apparent in
170 binary systems (logistic $f(x)$ with $\Delta = 0$, $x_{th} = 1$) and is not so apparent for systems with $f(x) =$
171 $\tanh(x)$. As might be expected, multistability for the logistic function with a steeper slope ($\Delta = 0.1$) is
172 more similar to that of the binary function, and the logistic function with the shallower slope ($\Delta = 0.2$) is
173 more similar to the tanh function.

174 While the impact of the response function, $f(x)$, on the results in Figure 3 suggest our findings
175 arise from more than a finite size effect, we wanted to test that possibility further. Therefore, we
176 assessed, for different networks with $s < 1$, how the probability of multistability depends on network
177 size. Our goal is to see how reliably cross-connections whose random strengths have a mean of zero
178 could, with increasing standard deviation, g , generate multistability that is absent with low g .
179 Simulation results of Figure 3 suggest a peak as a function of network size, N , in the likelihood of
180 reaching multiple final states from a fixed (large) number of initial conditions for some parameters.
181 However, without the exhaustive sampling of initial conditions, which becomes unfeasible at large N ,
182 our lack of multistability at large- N is not conclusive of its absence. To proceed further, we calculate

183 results for the infinite- N system in Section 3 and develop analysis of binary networks that allows for
184 finite- N approximations in Section 4.

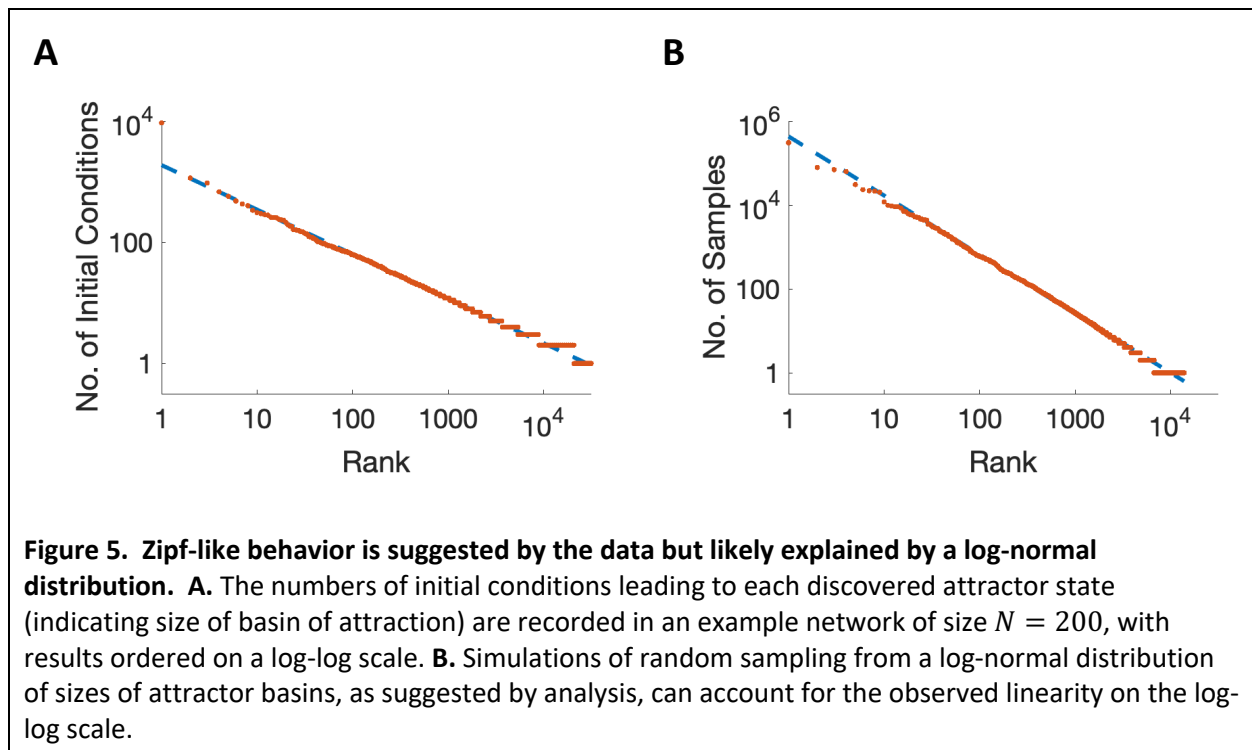


185

186 2.3. Distribution of size of basins of attraction

187 A major goal in our simulations of neural circuits was to assess the number of stable states they
188 contained as a marker of their information-carrying capacity. In systems with zero cross-connection and
189 strong enough self-interaction ($s > 1$) that each unit could be independently bistable, the number of

190 states is trivially 2^N , the maximal possible for our systems. However, without cross-connections, the
191 ability of such circuits to process information in a history dependent manner vanishes. Also, given the
192 unlikelihood that neural circuits operate in a regime with distinct bistable units, our focus was on circuits
193 with $s < 1$ such that individual units were not bistable, but with sufficiently strong $g > 0$ such that the
194 network could be multistable. For the results in this subsection, we focus on such systems with binary
195 units, $f(x) = \text{Heaviside}(x - x_{th})$, with $x_{th} = 1$.



196
197 For small circuits, $N < \sim 25$, one can sample initial conditions systematically with each of the 2^N
198 combinations of high/low activity per unit tested, but for larger networks one can sample only a subset
199 of all the possible initial conditions. We find that some states have vastly more initial conditions
200 reaching them compared to others. Indeed, if we take the number of randomly chosen initial conditions
201 that reach a particular fixed point as an indication of the size of the corresponding basin of attraction,
202 we find what appears to be a power law (Figure 5). Such a result suggests that in the absence of

203 exhaustive sampling, we will inevitably miss some of the smallest basins of attraction (note that the
204 frequency of visits ranges over 5 orders of magnitude in Figure 5A).

205 The suggestion of Zipf's Law in Figure 5A led us to consider theoretical reasons for producing
206 such a distribution. Our conclusion is that the apparent power-law is an artefact produced by sampling a
207 log-normal distribution with a very large width. Our reasoning is as follows. For any stable state some
208 units can have a level of input such that the unit could be stably active or inactive (assuming no change
209 in input from others). If N is large, the switching of such a unit does not strongly change the net input to
210 all other units, so another stable state is reached. Such a switch to a different state indicates the
211 crossing to a different basin of attraction. Across all of the distinct states in the network, the number of
212 units with inputs in the bistable range allowing for such stable switching is distributed as a Binomial (if
213 we ignore correlations), which is approximately a Normal distribution at large N . Or, equivalently, the
214 number of units that can be switched in any state *without* changing to a different basin of attraction
215 follows a Normal distribution across states. Additionally, the number of combinations of switching a unit
216 without producing a new stable state is approximately exponential in the number that can be
217 individually switched. Combining the two heuristics would suggest a log-normal distribution of sizes of
218 basins of attraction.

219 To test if our results in Figure 5A were compatible with a log-normal distribution, we generated
220 10^7 samples from a fictitious system with 10^7 states whose sizes, x , were distributed as $P(x) \propto$
221 $\exp\left[-\frac{[\ln(x)-\mu]^2}{2\sigma^2}\right]$, with $\mu = 30$ and $\sigma = 6$. Such a system was chosen to resemble the statistics, in terms
222 of numbers of states selected and maximum number of selections for any network of one of our random
223 systems with $N = 200$.

224 The results of our random sampling of a fictitious log-normal system in Figure 5B, indicate that
225 the observation of Zipf's Law from sampling basins of attraction is, indeed, compatible with a log-normal
226 distribution of attractor-basin sizes. The reason being that many of the small basins, whose expected

227 number of visits is less than one, either do not appear in the sampling at all, or appear once or twice,
228 and therefore increase the number of low-frequency states in a manner that “straightens out” the
229 inverted parabola that would be seen following an exhaustive sampling of the entire state space.

230 In summary, the observed frequency of visits of different attractor states follows an
231 approximate power law, but such behavior is most likely the consequence of sub-sampling of a
232 distribution which is approximately log-normal.

233

234 3. Mean-field theory

235 We followed the methods of others (Ahmadian et al., 2015; Stern et al., 2014) to develop a mean-field
236 theory for the large- N limit ($N \rightarrow \infty$) of each system. The following description of the method has a
237 slightly different emphasis from those of others, in part to connect to the alternative methods of section
238 4, and in part because our focus is on the existence of multiple stable fixed points rather than on the
239 more general dynamics of the system.

240 In the large- N limit, because each individual connection strength scales to zero, the impact of
241 small motifs (*e.g.*, small subsets of units with net positive interactions) and correlations in activity
242 between units becomes negligible. Therefore, the existence and stability of any state can be assessed by
243 assuming all units receive input sampled from the same distribution arising from the sum of connection
244 strengths multiplied by the activities of units. In small systems, the common scenario that units with
245 positive connections are more likely to be coactive together, renders the simplifying assumption
246 inaccurate. The large- N limit is also then (as stated in (Stern et al., 2014)) equivalent to averaging over
247 all realizations of the connectivity matrix, J_{ij} , which removes any correlation between individual units.

248 In the absence of any unit-specific identity, the unit label can be dropped from the formalism
249 and the dynamical mean field equation is one for the distribution of activations represented by the

250 variable $\mathbf{x}(t)$ in the face of a distribution of inputs given by a new variable, $\boldsymbol{\eta}(t)$, which we call the
 251 “field”:

$$\frac{d\mathbf{x}}{dt} = -\mathbf{x} + s\mathbf{f}(\mathbf{x}) + \boldsymbol{\eta}(t). \quad (2)$$

252 Self-consistency requires that the field, $\boldsymbol{\eta}(t)$, is produced by the sum of the product of distribution of
 253 activities, $\mathbf{f}(\mathbf{x}(t))$ (which result from the distribution of activation variable, $\mathbf{x}(t)$) multiplied by the
 254 connectivity matrix. Given the Central Limit Theorem, $\boldsymbol{\eta}$ is distributed as a Gaussian (a result justified
 255 more rigorously by others (Sompolinsky & Crisanti, 2018)) and given the lack of correlations between
 256 activity and connectivity in the large- N limit, we have:

$$\langle \boldsymbol{\eta} \rangle = \langle \mathbf{f}(\mathbf{x}) \rangle \langle gJ \rangle = 0 \quad (3)$$

257 and

$$\langle \boldsymbol{\eta}^2 \rangle = \langle \mathbf{f}^2(\mathbf{x}) \rangle \langle g^2 J^2 \rangle = g^2 \langle \mathbf{f}^2(\mathbf{x}) \rangle, \quad (4)$$

258 where we have used J to represent the $N \rightarrow \infty$ limit of $\frac{1}{\sqrt{N}} \sum_{j=1, j \neq i}^N J_{ij}$.

259 Fixed points of the dynamics (Eq. 2) arise for the distribution of activations, \mathbf{x} , where

$$\mathbf{x} - s\mathbf{f}(\mathbf{x}) = \boldsymbol{\eta}. \quad (5)$$

260 We define the variance of the zero-mean Gaussian distribution of $\boldsymbol{\eta}$ as σ^2 , such that

$$P(\boldsymbol{\eta}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\boldsymbol{\eta}^2}{2\sigma^2}\right), \quad (6)$$

261 where σ^2 must be calculated self-consistently from Eq. 4. For the system to possess multiple attractor
 262 states, the above set of equations (2)-(6) must have multiple solutions, and those solutions must
 263 correspond to stable states.

264 Multiple solutions arise from Eq. (5) if the function $\mathbf{x} - s\mathbf{f}(\mathbf{x})$ is non-monotonic. Given the
 265 neural response function, $\mathbf{f}(\mathbf{x})$, has zero slope at very negative or very positive \mathbf{x} , the function $\mathbf{x} -$
 266 $s\mathbf{f}(\mathbf{x})$ has slope of +1 at these extremes and is therefore non-monotonic if for any value of \mathbf{x} we have
 267 $\mathbf{f}'(\mathbf{x}) > 1/s$. That is, multistability is possible if $s > \max \mathbf{f}'(\mathbf{x})$. Hence the result in (Stern et al., 2014)

268 that if $f(x) = \tanh(x)$ with a maximum gradient of 1, multistability is only possible if $s > 1$. For $f(x) =$
269 $\frac{1}{1+e^{\frac{x_{th}-x}{\Delta}}}$ the requirement is $s > 4\Delta$.

270 It is important to note that multiple self-consistent solutions of Eq. (4) for the variance of the
271 field, η , are also possible. Indeed, for the logistic input-output function, a solution with low variance
272 corresponding to a quiescent, or low-activity state (in which activities of units are tightly clustered
273 around $f(0)$) can coexist with a solution of greater input-variance. We assess both the stability of the
274 solution with minimal activation (and therefore minimal variance of the field) as well as the existence of
275 and stability of distinct solutions with higher activation when determining which states exist for a given
276 set of parameters (see Appendix 3 for methods).

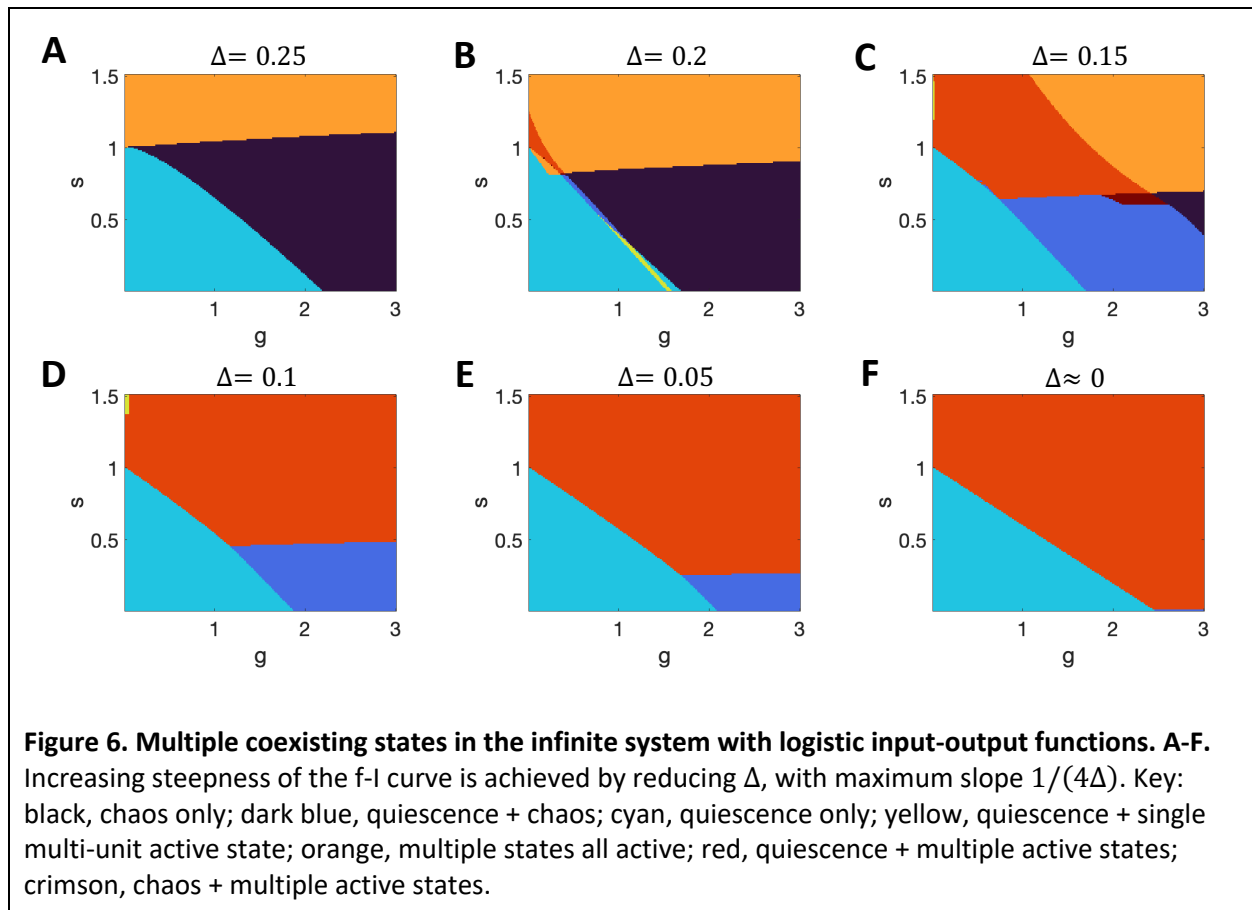
277

278 3.1. Phase diagram for networks with logistic single-unit response functions

279 In Figure 6, we show how the phase diagram depends on the slope of the input-output function, $f(x) =$
280 $\frac{1}{1+e^{\frac{x_{th}-x}{\Delta}}}$, with the panels from A to F depicting results of increasing steepness (by lowering Δ). The final
281 panel ($\Delta = 0$) is produced by methods described in the next section. In all cases x_{th} is adjusted according
282 to Eq. A1, to ensure single-unit bistability at $s = 1$. As can be seen, the minimum level of s allowing for
283 multiple stable active states falls in proportion to Δ , and in the range $4\Delta < s < 1$, the system is
284 quiescent at $g = 0$, but with increasing g becomes multistable. In all cases, and for all values of s , the
285 expected transition to chaos arises with large enough g , though that transition is not always visible in
286 the parameter ranges shown.

287 In systems with $0 < \Delta < 0.25$ (Figure 6 B-E) we find ranges of parameters for which the field, η ,
288 does indeed have two self-consistent solutions. Most commonly, at low s , the cyan regions indicate the
289 presence of a stable quiescent state with an unstable active state, that is the coexistence of inactivity

290 and chaos in a given network. In a smaller range of parameters, the yellow region (Figure 6B) indicates
 291 the coexistence of a stable quiescent state with a stable active state. Such multistability exists even in



292 the absence of cross-connections ($s = 0$) and concurs with our simulation results in Section 2.1. Indeed,
 293 the region of multistability spans the value of $g = 1.55$, observed at larger- N in simulations. Therefore,
 294 even without the self-excitation needed for individual units to be bistable with sufficient input, the
 295 network can possess multiple stable states, with the two distinct states resulting from and causing two
 296 distinct population-mean (and mean-square) firing rates and two distinct population input distributions.

297 298 3.1.1 Accounting for extreme tails of a Gaussian in the Infinite System

299 Our definition of the quiescent state contains a requirement that all units have activity of less than half
 300 of their maximum. In practice, in systems where bistability is possible ($4\Delta < s < 1$) the quiescent state

301 requires that all units are stable on the lower branch of the bifurcation curve. In the infinite system, any
302 requirement of *all* units raises a subtle issue that we address in this subsection (and in Appendix C and
303 Supplementary Figures 3 and 5).

304 The field, η , which indicates the probability of any unit receiving a given input, follows a
305 Gaussian distribution. When all firing rates are very low, the variance, σ^2 , of the Gaussian distribution
306 for η is very low but is non-zero. The probability of a unit receiving input with a magnitude $Z\sigma$, that is
307 many times, Z , greater than the standard deviation, σ , is vanishingly small (*e. g.* if $Z = 6$ the probability
308 is less than 10^{-9} and if $Z = 9$ the probability is less than 10^{-18}). However, for any finite Z the
309 probability is strictly non-zero for any finite-level of input, so an infinite system will always have units
310 whose inputs exceed that value. Therefore, in a system in which $s > 4\Delta$, the quiescent state is unstable
311 for an infinite system, unless $g = 0$ precisely. Yet, for any biologically feasible circuit we can define a
312 Z_{max} and require the bifurcation points, η^* , to be within the range $-Z_{max}\sigma \leq \eta^* \leq Z_{max}\sigma$, in order for
313 the quiescent state to be defined as unstable. In this manner, we can study a system in which we have
314 ignored correlations (an approach strictly only correct in the infinite- N limit) but at the same time define
315 states that would be present in a large finite system with results accurate (to 1 part in 1000) for sizes up
316 to $N = 10^6$ (with $Z_{max} = 6$) or even $N = 10^{15}$ (with $Z_{max} = 9$).

317 A similar issue arises when we consider whether a system has multiple stable active states. The
318 number of such states depends (exponentially) on the number of units receiving input between the two
319 bifurcation points of $\mathbf{x}(\eta)$ such that the unit could be either active or inactive for that level of input.
320 Whenever there is a pair of bifurcation points ($4\Delta < s < 1$), the argument from the previous paragraph
321 again indicates that in an infinite system there is always a unit with input in that range. However, while
322 in the infinite system the network is multistable, for any realistic system—even a large one—it may be
323 very unlikely that any unit receives sufficiently extreme input, so we use the same value of Z_{max} to
324 indicate multistability as we do for stability of the quiescent state.

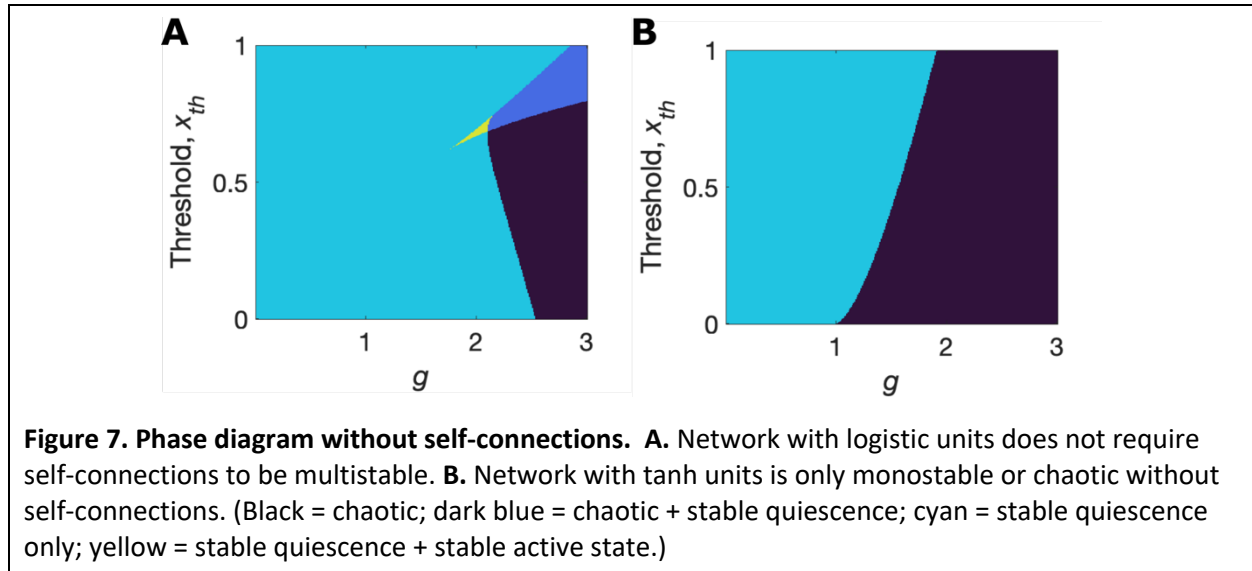
325 Therefore, in Supplementary Figure 3, we replot the phase diagrams for distinct levels of Z_{max} ,
326 while using a standard of $Z_{max} = 6$ for most phase diagrams. For example, with $Z_{max} = 6$, a large
327 network of 10^6 units would only have a probability of 0.001 of behaving differently from that indicated
328 in the phase diagram (or a network of 10^3 units would have a 1 in a million chance of behaving
329 differently—the fewer the units in a network, the less likely at least one of those units receives
330 excessively high input). In the limit of $Z_{max} \rightarrow \infty$, the system has a discontinuity moving away from the
331 y-axis, as even while the region of multistability approaches $g = 0$ (for $4\Delta < s < 1$) throughout this
332 paper we have set the threshold, x_{th} , such that disconnected units are only bistable if $s > 1$.

333 By contrast, when analyzing the network with tanh units of low Δ and higher x_{th} —that is a
334 steeper f-I curve, but with single-unit bifurcation maintained at $s = 1$ —the output of units with zero
335 input is close to -1, rather than 0. Such maximally negative output produces a larger variance of inputs
336 across units, such that multistability is more common at very low cross-connection strength
337 (Supplementary Figure 4). Therefore, while the choice of Z_{max} still impacts the phase diagram for those
338 reasons discussed above, it does so to a much smaller extent for tanh units (Supplementary Figure 5).

339

340 3.2. Multistability without self-connections

341 Multistability without self-connections (*i.e.*, with $s = 0$) is present in all networks with logistic response
342 functions if we allow x_{th} to vary (or equivalently apply uniform input). To demonstrate this, in Figure 7A
343 we show the phase diagram as a function of x_{th} and g for a system with $s = 0$ and $\Delta = 0.25$ —in this
344 case $f(x) = \frac{1}{2}[1 + \tanh(x)]$. Systems with different Δ produce identical figures if the two axes are
345 scaled by the same factor as Δ . As can be seen, the region of coexistence of quiescence with chaos is
346 contiguous with and extends a region of coexistence of quiescence with an active stable state. These
347 two regions, which depend on multiple stable solutions for the self-consistency of the field, $\boldsymbol{\eta}$, are not
348 present if the response function is $f(x) = \tanh(x)$ (Figure 7B).



349

350 The two response functions, $f(x) = \frac{1}{1+e^{-\frac{x-x_{th}}{\Delta}}}$ and $f(x) = \tanh\left(\frac{x-x_{th}}{\Delta}\right)$, have a key difference

351 that leads to them producing qualitatively different behavior. For the logistic function, the minimal
 352 absolute value of $f(x)$ coincides with the minimal gradient of the function (if $|f(x)|$ is low then $f'(x)$ is
 353 low) whereas for the tanh function the opposite is true (if $|f(x)|$ is low then $f'(x)$ is near its maximum).
 354 Hence for the logistic function, it is possible for a narrow range of inputs to produce a stable narrow set
 355 of low firing rates (maintaining low inputs) while a solution with a large range of inputs leads to some
 356 much larger stable firing rates (maintaining high inputs) given the supralinearity of the response
 357 function. However, for the tanh function, the marginal feedback decreases with a change in rate from
 358 zero, so only one solution can exist.

359

360 4. Analysis of networks of units with binary response functions

361 For a system with binary units, $f(x) = Heaviside(x - x_{th})$, the analysis simplifies, because a state is
 362 stable if all active units have input from other active units exceeding $x_{th} - s$ and all inactive units have
 363 summed input from the active units less than x_{th} . We define k as the number of active units, each with

364 an activity of 1, so the above requirements on network inputs correspond to the sum of $k - 1$ of the
 365 connection strengths to each of k active units and to the sum of k connection strengths to each of the
 366 $N - k$ inactive units.

367 Our main approximation is to treat these sums of connections strengths as independent draws
 368 from a Gaussian distribution with mean of zero and whose variance is $\frac{g^2(k-1)}{N}$ and $\frac{g^2k}{N}$ respectively. We
 369 then assume a solution with k active units exists if, given N independent draws from a Gaussian of unit
 370 variance, the top k draws, when multiplied by $\sqrt{\frac{g^2(k-1)}{N}}$ are greater than $x_{th} - s$ (high input to active
 371 units) while the remaining $N - k$ draws, when multiplied by $\sqrt{\frac{g^2k}{N}}$ are less than x_{th} . This is equivalent to
 372 the requirement that the $(k + 1)_{th}$ greatest sample out of N samples, $X_{k+1,N}$, from a unit-variance,
 373 zero-mean Gaussian distribution lies in the range:

$$(x_{th} - s) \sqrt{\frac{N}{g^2(k-1)}} < X_{k+1,N} < x_{th} \sqrt{\frac{N}{g^2k}} \quad (7)$$

374 where we have assumed $x_{th} - s > 0$ and $x_{th} > 0$ (which holds in our standard system with $x_{th} = 1$ so
 375 long as $s < 1$, and which is the parameter region in which we have greatest interest).

376 Such a requirement can be calculated using the methods of order statistics (David & Nagaraja,
 377 2003), which we follow for the Gaussian distribution, defining

$$378 \quad P(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

$$379 \quad P_+(x) = \int_x^\infty P(x') dx' = \frac{1}{2} \operatorname{erfc}\left(\frac{x}{\sqrt{2}}\right)$$

$$380 \quad P_-(x) = \int_{-\infty}^x P(x') dx' = 1 - P_+(x) = \frac{1}{2} + \frac{1}{2} \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right).$$

381 such that

$$P(X_{k+1,N} = x) = N \binom{N-1}{k} P(x) [P_+(x)]^k [P_-(x)]^{N-k-1}. \quad (8)$$

382 Therefore, we have a stable state with k of N units active with a probability $P(k, N)$ given by:

$$P(k, N) = N \binom{N-1}{k} \int_{(x_{th-s})\sqrt{\frac{N}{g^2(k-1)}}}^{x_{th}\sqrt{\frac{N}{g^2k}}} P(x)[P_+(x)]^k [P_-(x)]^{N-k-1} dx. \quad (9)$$

383 We then calculate the probability a system has at least one stable state with multiple active units, and

384 so is multistable (as the quiescent state is always stable for $x_{th} > 0$) for a given system size, N , via:

$$P(N) = 1 - \prod_{k=1}^N [1 - P(k, N)], \quad (10)$$

385 (there is only an absence of multistability if it is absent for all possible k). Again, (10) is an approximation

386 as it assumes $P(k, N)$ is independent for different values of k . We will see that in the large- N limit the

387 approximation becomes exact, as $P(k, N)$ becomes either 0 or 1, so that $P(N)$ is also either 0 or 1, and

388 we have multistability with probability 1 if and only if it arises with probability 1 for some value of k ,

389 that is $P(N) = \max P(k, N)$.

390

391 4.1. Finite- N results of analysis with binary units

392 Our simulation results presented in Section 2 (Figure 4) suggest that, for some parameters, networks of

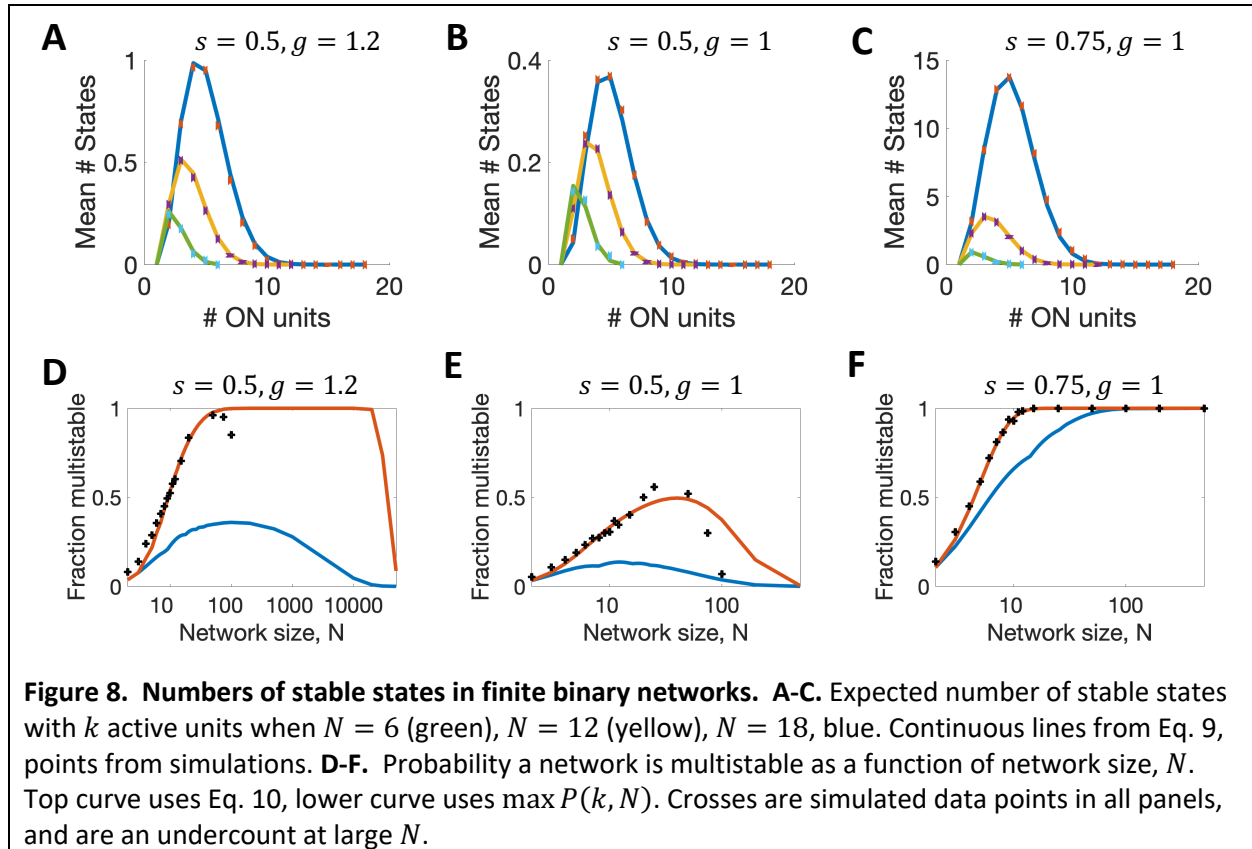
393 intermediate size have the greatest probability of multistability. Given the increasing likelihood of

394 missing stable states as N increases using simulation methods, that simulation result may be incorrect

395 due to undersampling at large N . Therefore, we use our approximate analytical methods for networks

396 with binary units to address the N -dependence of the probability of multistability, by solving Equations

397 (8)-(10) above.



398

399

400

401

402

403

404

405

406

407

408

In Figure 8A-C, we show that for small networks of $N = 6$, $N = 12$, and $N = 18$, for which we can exhaustively test all initial conditions and therefore find all stable states in simulations, the approximate analytical method (solid line, which plots Eq. 9) compares well with the simulated data (crosses). Moreover, in Figure 8D-F, when we use Eq. 10 (red lines) to estimate the probability the network is multistable, the simulated results (crosses) are remarkably close to the analytic approximation. Such a result is surprising, as one would expect a positive correlation across networks and the numbers of stable states. The blue lines in Figure 8D-F are the results for a correlation of +1, in which the network's probability of multistability is simply the maximum across possible states and is much farther from the data than the analysis assuming zero correlation (the red line). Nevertheless, across all methods, in Figure 8D-E, we do indeed find that the probability of multistability peaks at

409 intermediate network size, remarkably reaching values of approximately 1 for $s = 0.5$, $g = 1.2$, before
 410 falling to zero at large network size ($N > 30000$).

411

412 4.2. Large- N limit of system with binary units

413 In the large- N limit, the above equations (8)-(9) can be simplified. First, we note that

$$\binom{N-1}{k} [P_+(x)]^k [P_-(x)]^{N-k-1} = \binom{N-1}{k} [P_+(x)]^k [1 - P_+(x)]^{(N-1)-k} \quad (11)$$

414 is the Binomial probability for achieving k outcomes from $N - 1$ independent selections, with individual
 415 probability of outcome, $P_+(x)$. In general, the probability has a peak at the integer value of k closest to

416 $(N - 1)P_+(x)$, with a standard deviation of $\sqrt{(N - 1)P_+(x)[1 - P_+(x)]}$.

417 If we define $f = \frac{k}{N} \cong \frac{k-1}{N}$ with $k \gg 1$ as well as $N \gg 1$, then the integration limits for $P(k, N)$ become

418 $\frac{x_{th}-s}{g} \sqrt{\frac{1}{f}}$ and $\frac{x_{th}}{g} \sqrt{\frac{1}{f}}$. Also, for large N , the Binomial probability term approaches a Dirac delta-function

419 at the value of $f = P_+(x)$, as it's standard deviation in f scales as $1/\sqrt{N}$.

420 Therefore, in the large- N limit, $P(f) = 1$ if the integration range over x contains the value

421 where $f = P_+(x)$ and $P(f) = 0$ otherwise. Algebraically this becomes a requirement that f lies

422 between two thresholds, Θ_1 and Θ_2 , which each depend on f :

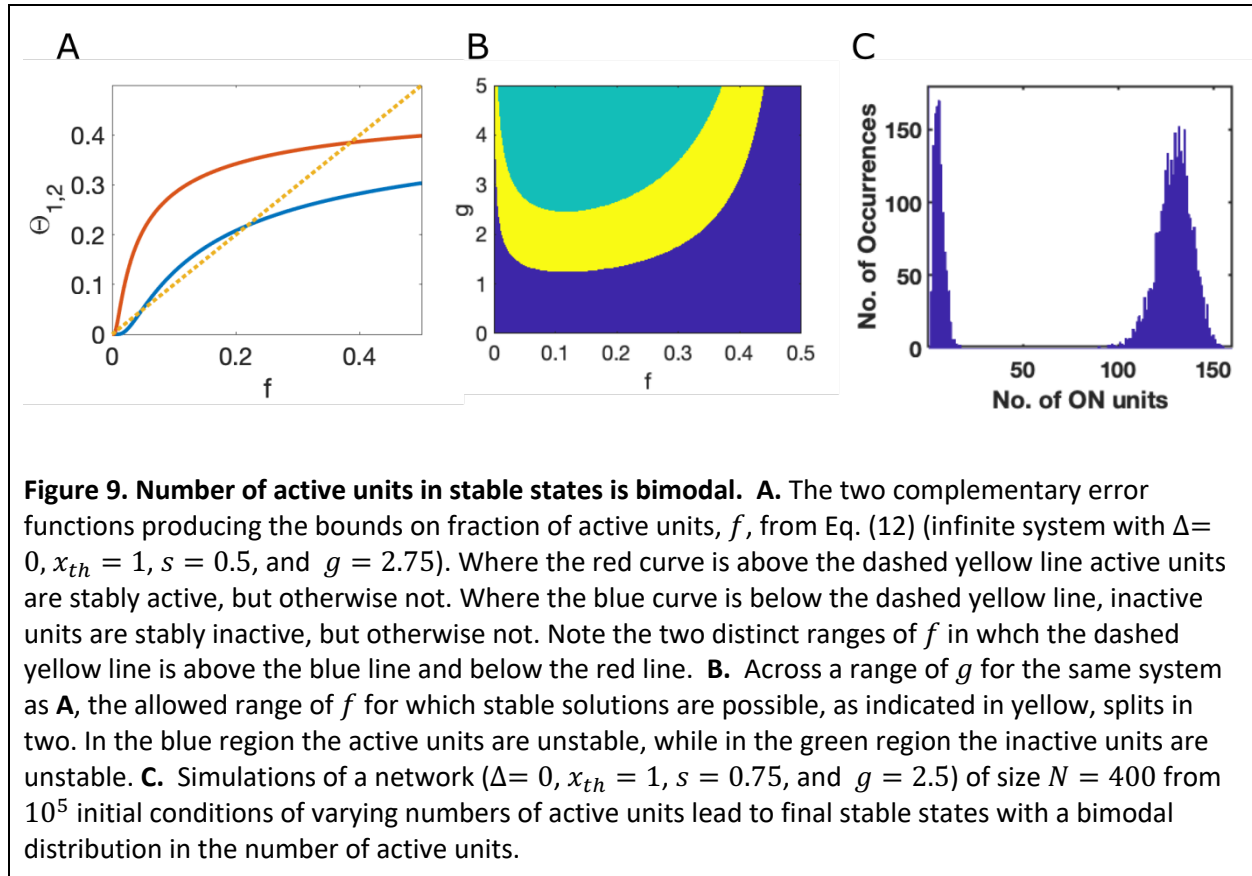
$$P(f) = 1 \text{ if } \frac{1}{2} \operatorname{erfc}\left(\frac{x_{th}}{g\sqrt{2f}}\right) = \Theta_1 < f < \Theta_2 = \frac{1}{2} \operatorname{erfc}\left(\frac{x_{th}-s}{g\sqrt{2f}}\right). \quad (12)$$

423 Figure 9A indicates the region of these inequalities as a function of f for the specific values of $x_{th} = 1$,

424 $s = 0.5$, and $g = 2.75$, with Figure 9B showing that for a wide range of g the possible numbers of

425 active units in a stable state splits into two distinct ranges. Simulation results in Figure 9C demonstrate

426 such bimodality in numbers of active units for a similar network ($s = 0.75$ and $g = 2.5$) with $N = 400$.



427 Given that at $f = \frac{1}{2}$ it is always true that $\frac{1}{2} \operatorname{erfc}\left(\frac{x_{th}}{g\sqrt{2f}}\right) < f$ (because of the limited range of the
428 complementary error function) the criterion for Eq. 12 to have a solution for some value of f
429 is the requirement $f < \frac{1}{2} \operatorname{erfc}\left(\frac{x_{th}-s}{g\sqrt{2f}}\right)$ for some f , which leads to a minimum value of $g = g^*$ at which
430 the lines $y = x$ and $y = \frac{1}{2} \operatorname{erfc}\left(\frac{x_{th}-s}{g\sqrt{2x}}\right)$ meet at a tangent. The critical value occurs where $\left(\frac{x_{th}-s}{g} \cong \frac{1}{2.457}\right)$
431 such that, as a function of s and with $x_{th} = 1$, multistability arises in this system if $g > g^*(s) \cong$
432 $2.457(1 - s)$ (Figure 6F).

433 Notice that as s approaches zero, the range of f with allowed solutions shrinks toward the line
434 where $f = \frac{1}{2} \operatorname{erfc}\left(\frac{x_{th}}{g\sqrt{2f}}\right)$. For $g > g^*$, there are two such solutions, which are distinct crossings of the
435 line. The distinct solutions indicate two separate ranges for the possible number of active units in stable
436 attractors. In the example shown in Figure 9B, solutions are possible in two ranges, either with a very

437 low fraction (<5%) of units active or with a fraction in the range of 30%-40% of units active. The two
438 distinct ranges are also visible from simulations as a bimodal distribution in the numbers of active units
439 in stable states following random initial conditions (Figure 9C).

440 A lower bound on the fraction of active units arises, because with few active units there is too
441 little network input to activate those units. The upper bound arises because half of the units receive net
442 negative input, so cannot be stably active if $s < 1$, and only a subset of those units receiving net positive
443 input receive an amount greater than $1 - s$, as needed to be stably active. The bounded region in which
444 active units have sufficient input to remain stably active can contain within it a separate bounded region
445 of instability (Figure 9) because the random network input can be sufficiently strong that some of the
446 inactive units (of which there are more than there are active units) receive too much input to remain
447 “off”.

448

449 5. Discussion

450 Firing rate models of neurons are valuable because they represent the likely states of a neural circuit in a
451 relatively simple manner and can be solved rapidly. The foundation of a firing rate model is the input-
452 output function of a neuron, which is typically designed to have bounded outputs over the domain of
453 inputs. For its ease of mathematical manipulation, the hyperbolic tangent function, $f(x) = \tanh(x)$, has
454 been used with great success, most notably for first demonstrating the transition from quiescence to
455 chaos as the strength of random cross-connections increases (Sompolinsky et al., 1988; Stern et al.,
456 2014). The negative portion of $\tanh(x)$, while it cannot correspond to negative firing rates, could be
457 considered representative of a group of mixed excitatory-inhibitory neurons in which the mean rate of
458 inhibitory neurons exceeds that of excitatory neurons.

459 Given the function $f(x) = \tanh(x)$ is simply a translated version of the function $f(x) =$
460 $\tanh(x - x_{th}) + 1$, one might expect that analysis of a system with units responding via the one
461 function would provide all the qualitative insight necessary to understand the behavior of a system with
462 units responding via the other function. However, this is not the case. A disconnect between the
463 behavior of a system of neurons with $f(x) = \tanh(x)$ and that of a system with $f(x) = \tanh(x) + 1$
464 has been shown by others (Figure 4b of (Touboul & Ermentrout, 2011)) whereby a Hopf bifurcation
465 disappears as the input-output function of neurons is parametrically shifted up toward non-negative
466 values. In our analyses, we find two qualitative changes. The first is a shift in phase boundaries leading
467 to the result that random cross-connections, whose mean value is zero, can produce multistability in a
468 system in which single units are not in of themselves bistable.

469 Second, we find the possibility of bistability via distinct stable solutions for the self-consistency
470 of the field. Alternative self-consistent solutions of the field can lead to multistability arising from
471 random, zero-mean, cross-connections even in systems without self-connections (Figure 3 and Figure
472 6B). The distinct self-consistent field solutions, with different variances in the input currents, correspond
473 to states with distinctly different numbers of active units. Figure 9 indicates a similar bimodality in the
474 numbers of active units in simulated binary-unit systems and is coupled with an analysis of how such
475 bimodality arises in the system.

476 We find a subtlety when taking the infinite limit of our system using the logistic input-output
477 function, with a strict discontinuity between results with $g = 0$ and those with $g = \epsilon$ (where g scales
478 the strength of cross-connections and ϵ is an infinitesimal positive quantity). The reason being that for
479 non-zero g , there is a non-zero (even if miniscule) probability that the within-circuit input to a unit,
480 which is drawn from a Gaussian with width proportional to g , is sufficiently strong to render that unit
481 bistable. However, when the bifurcation point is many tens of standard deviations above the zero mean
482 of the Gaussian distribution, the probability becomes infinitesimal and is irrelevant in any real or

483 simulated system, even with billions of units. For similar reasons, the strict mathematical limit has a
484 discontinuity when altering the width of the logistic function from $\Delta = 0$ to $\Delta = \epsilon$. If, instead of producing
485 a phase diagram, with a sharp boundary for multistability, we focused on the entropy of the system (the
486 log of the number of stable states) scaled by system size, N , such discontinuities would disappear as the
487 entropy would reduce continuously and smoothly (and rapidly) from the boundaries of multistability
488 shown in Figure 6, to a tiny value before becoming strictly zero at $g = 0$ or $\Delta = 0$.

489 Multistability, when exhibited as a set of discrete stable fixed points, may seem unlikely in any
490 cortical circuit given that activity is never static *in vivo*. However, a network based on multiple fixed
491 points, but with randomly timed transitions between them, can match the observed data in a number of
492 systems (Ballintyn et al., 2019; Ksander et al., 2021; La Camera et al., 2019; Mazzucato et al., 2019;
493 Miller, 2016; Miller & Katz, 2010; Moreno-Bote et al., 2007; Recanatesis et al., 2022). Moreover,
494 analyses of patterns of neural spiking *in vivo* have, in many cases, shown that a discrete state-based
495 formalism better matches the data than a formalism assuming continuously changing, graded activity
496 (Abeles et al., 1995; Miller & Katz, 2010, 2011; Ponce-Alvarez et al., 2012; Sadacca et al., 2016;
497 Seidemann et al., 1996).

498 While the strengths of connections between units are treated as independent random variables
499 for ease of analysis in this paper, in practice there is internal structure in the connectivity among
500 neurons, even between excitatory pyramidal cells (Song et al., 2005; Stepanyants & Chklovskii, 2005).
501 Moreover, connections from cortical neurons typically have fixed sign (all excitatory or all inhibitory)
502 according to neuron class, a feature that can change the behavior of random networks (Rajan & Abbott,
503 2006). In our work, we consider a firing rate model unit as representing the mean rate of a cluster of
504 many neurons (as is necessary to omit the pulsatile spike interaction from simulations) so the net
505 interaction between units can be of either sign according to whether the dominant connections are
506 excitatory-to-excitatory, or excitatory-to-inhibitory, etc. Moreover, much of the nonrandom cortical

507 structure can be accounted for by considering the intra-cluster connectivity to be distinct from the inter-
508 cluster connectivity (Bourjaily & Miller, 2011) as we do here.

509 Our main conclusion is that multistability can be produced via random, zero-mean cross-
510 connections in neural circuits without the exceptionally strong self-connections needed to produce
511 bistability in a single cluster of neurons (a unit in a firing-rate model) so long as the neurons without
512 input have a low firing rate and if rate increases supralinearly with low input.

513 Code Availability

514 MATLAB codes used to produce the results in this paper are available for public download at
515 <https://github.com/primon23/Multistability-Paper>.

516 Acknowledgments

517 The authors are grateful to NIH-NINDS for support of this work via R01 NS104818 and to the Swartz
518 Foundation for a fellowship to SQ. SM is grateful to Merav Stern for helpful conversations in the early
519 stages of this work.

520 Appendix 1: Monte Carlo simulation method

521 Our standard procedure was to simulate 100 different realizations of the connectivity matrix to produce
522 100 random networks for a given parameter combination. For each connectivity matrix, we then
523 completed sets of multiple trials, each trial with a distinct initial condition (100 trials for perturbation
524 analysis in Figure 2 and for scaling g in Figure 3; 200 trials for parameter grids in Figure 4; and 10^6 or 10^5
525 trials respectively for the networks with binary units in Figures 8 and 9) . For the small ($N \leq 25$)
526 networks with binary units in Figure 8 all 2^N combinations of initial conditions were used with each unit
527 at an initial rate of its minimum or maximum.

528 The continuous models were simulated using MATLAB's ode45 function. Each trial was
529 simulated until either a maximum simulation time was reached (5,000 τ for Figure 3 and 10,000 τ for
530 Figure 4), or until a stopping condition was reached in the case that the maximum \dot{x}_i at a give timestep
531 was less than $2 * 10^{-6}$. If this stopping condition was reached, then the activity was considered to have
532 reached a stable state because the network possessed a point attractor at that set of firing rates.
533 Logistic units were classified as active if their firing rate exceeded 0.5. Tanh units were considered active
534 if the absolute value of their rate exceeded 0.001. For the continuous models, typically the first trial was
535 initialized with inputs near zero, to test if the quiescent state was stable. For all subsequent trials, the
536 initial rates of the units were set to a uniform random distribution over 0 to 1 and transformed by a
537 logistic function with $x_{th} = 0.5$ and $\Delta = 0.1$.

538 For the perturbation analysis, each trial of each network was simulated for a full 21,000 τ . Then,
539 at each of 100 linearly spaced time points between 20,000 and 20,800 τ 10% of the units' firing rates
540 were randomly perturbed upwards or downwards by 10^{-5} and the simulation was then continued from
541 each such perturbed state for 200 τ . The root mean squared (RMS) deviation of the perturbed
542 simulation from the original simulation quantified the extent to which the perturbation caused a
543 divergence in activity. The median RMS deviation over the 100 perturbations was then used to classify
544 each trial as a point attractor, a limit cycle, or chaotic. The median RMS deviation exponentially decayed
545 for point attractors, exponentially increased for chaos, and increased but reached a plateau at a low
546 level for limit cycles. Classification thresholds were set based on the R^2 of a linear fit to the exponential
547 RMS deviation and the magnitude of the RMS deviation averaged between 190 to 200 τ post-
548 perturbation. Trials with final RMS deviations below half the magnitude of the initial perturbation and
549 with no units having a change in their firing rate exceeding 10^{-4} in the last 10 τ of the unperturbed
550 simulation were classified as point attractors. To classify trials as chaotic vs limit cycles, a classification
551 boundary was determined as a function of each trials' linear fit R^2 and final RMS deviation. Trials above

552 the line $RMS\ deviation = 0.025 * e^{(-0.125 * R^2)}$ were classified as chaotic. This boundary allows the
 553 separation between these two dynamics because it accounted for both chaotic trials that very quickly
 554 converged to a large RMS deviation (large RMS deviation and low R^2) and chaotic trials that had a slower
 555 exponential increase in their RMS deviation (lower RMS deviation at 190 to 200 τ but high R^2). Final
 556 activity states of the unperturbed simulations were used to confirm these classifications.
 557

558 Appendix 2: Choice of single-unit input threshold

559 For comparison across systems with distinct single-unit input-output functions, $f(x)$, we adjust the
 560 offset, x_{th} , such that a single unit becomes bistable with self-connection strength of $s = 1$, in all cases.

561 For the logistic function, such a requirement means that a saddle-node bifurcation occurs at $s =$
 562 1 , with unstable and stable fixed points colliding at x^* given by $-x^* + sf(x^*) = 0$ such that $x^* = f(x^*)$
 563 and $\frac{d}{dx}[-x + sf(x)]_{x^*} = 0$ such that $\frac{df(x)}{dx}\Big|_{x^*} = 1$. Combining these equations and using the result for
 564 the logistic function that $\frac{df(x)}{dx} = \frac{1}{\Delta} f(x)[1 - f(x)]$ leads to the requirement:

$$x_{th} = \frac{1}{2} + \sqrt{\frac{1}{4} - \Delta} + \Delta \ln \Delta - 2\Delta \ln \left(\frac{1}{2} + \sqrt{\frac{1}{4} - \Delta} \right). \quad (A1)$$

565 For the binary response function, $f(x) = Heaviside(x - x_{th})$, we have $x_{th} = 1$, which can be seen
 566 from the above equation in the limit $\Delta \rightarrow 0$.

567 For the hyperbolic tangent function, $f(x) = \tanh\left(\frac{x - x_{th}}{\Delta}\right)$, a similar derivation leads to

$$x_{th} = \sqrt{1 - \Delta} - \frac{\Delta}{2} \ln \left(\frac{1 + \sqrt{1 - \Delta}}{1 - \sqrt{1 - \Delta}} \right), \quad (A2)$$

568 which yields $x_{th} = 0$ if $\Delta = 1$, matching the simplest anti-symmetric response function, $f(x) = \tanh(x)$,
 569 and as with the binary response function, $x_{th} = 1$ if $\Delta = 0$.

570

571 Appendix 3: General mean-field methods

572 To test whether a distribution of the interacting variables, \mathbf{x} , produces a stable fixed point, it is
 573 necessary to obtain information about the eigenvalues of the Jacobian matrix of the dynamical
 574 equations expanded linearly about the fixed point (Strogatz, 2015). If all such eigenvalues have a
 575 negative real part then the fixed point is stable. Linearization around a fixed point, \mathbf{x}^* , yields

$$\frac{d\mathbf{x}}{dt} = -\mathbf{x} + s\mathbf{D}\mathbf{x} + \frac{g}{\sqrt{N}}\mathbf{J}\mathbf{D}\mathbf{x} \quad (\text{A3})$$

576 where \mathbf{D} is a diagonal matrix with elements equal to the corresponding derivatives of the input-output
 577 function, $f'(\mathbf{x}^*)$, and \mathbf{J} is the unit variance, zero mean, Gaussian connectivity matrix.

578 We follow the methods of others (Ahmadian et al., 2015; Stern et al., 2014) who showed that
 579 eigenvalues of such a system are found at the complex values, z , where

$$\text{Tr} \left[(M_z M_z^\dagger)^{-1} \right] \geq 1$$

581 with

$$M_z = \frac{(z + 1 - W_S^{EE} f'(\mathbf{x}^*))}{g f'(\mathbf{x}^*)}.$$

583 In the large- N limit the sum within the Trace become an integral over the distribution of activations, \mathbf{x} ,
 584 to yield the criterion (Ahmadian et al., 2015; Stern et al., 2014):

$$\int dx P(x) \frac{g^2 [f'(x)]^2}{|z + 1 - s f'(x)|^2} = \int d\eta P(\eta) \frac{g^2 [f'(x(\eta))]^2}{|z + 1 - s f'(x(\eta))|^2} \geq 1. \quad (\text{A4})$$

585 As noted by (Stern et al., 2014), for the system to be stable we require that Equation A4 is not satisfied
 586 for any z with $\text{Re}[z] > 0$, which allows us to assess the case where $\text{Re}[z] = 0$ and note that any non-
 587 zero contribution to $\text{Im}[z]$ increases the absolute value of the denominator above, so if there are no
 588 eigenvalues with $z = 0$ there cannot be any on the imaginary axis. Therefore, in general we require, for
 589 there to be no eigenvalues with positive real part that

$$\frac{1}{\sqrt{2\pi\sigma^2}} \int d\eta \exp\left(-\frac{\eta^2}{2\sigma^2}\right) \frac{g^2 [f'(x(\eta))]^2}{[1 - sf'(x(\eta))]^2} < 1, \quad (\text{A5})$$

590 where we have substituted for $P(\eta)$ and $\sigma^2 = g^2 \langle f^2(x) \rangle$. We have also assumed that the function in
 591 the denominator, $1 - sf'(x(\eta))$, is positive, as any negative portion of the function means there is a
 592 divergent positive contribution to the integral for some z with $Re[z] = sf'(x(\eta)) - 1 > 0$.

593 We are interested in cases of multistability, where the activations, $x(\eta)$, can have more than
 594 one value based on the solutions $x - sf(x) = \eta$ for some values of η . This requires that $x - sf(x)$ is a
 595 non-monotonic function, which occurs if $\max [f'(x)] > 1/s$ (to produce a region of negative slope in
 596 the function $x - sf(x)$). The need for a region of negative slope arises because in all cases considered
 597 here at large positive or negative values of x , $f'(x) = 0$ and $x - sf(x)$ has a slope of +1. In cases of
 598 multiple solutions for $x(\eta)$, care must be taken in the choice of $x(\eta)$, as while stability is enhanced by
 599 choosing the solution with the lower value of $f'(x(\eta))$, such a choice can lead to the lower value of
 600 $f^2(x)$ for some input-output functions (but not if $f(x) = \tanh(x)$) which can lead to the self-consistent
 601 solution for the distribution of η to become too narrow to support multistability, as discussed below.

602 In the logistic networks, $f(x) = \frac{1}{1 + \exp\left(\frac{x_{th} - x}{\Delta}\right)} \approx \exp\left(\frac{x - x_{th}}{\Delta}\right)$ for $x \ll x_{th}$, is never exactly zero.

603 Therefore the Gaussian distribution of η will always have non-zero variance for $g > 0$ and, even if the
 604 distribution is narrow with very small variance, the distribution always retains some vanishingly small
 605 but non-zero density at the values of η required to support multiple solutions of $x(\eta)$ if $s > 4\Delta$.
 606 However, if bifurcation points in $x(\eta)$ require levels of the Gaussian-distributed η that are many
 607 standard deviations from its mean of zero, such solutions give exponentially small probability of
 608 multistability in a finite network, so are unlikely to be observed in practice. Therefore, we set a
 609 threshold, $Z_{max}\sigma$, in terms of the number, Z_{max} , of standard deviations, σ , of the field of inputs, η , such
 610 that if both bifurcation points, η^* , are beyond the threshold ($\eta^* < -Z_{max}$ or $\eta^* > Z_{max}$) we ignore
 611 both the extra solutions and any instability they cause. To clarify the result of such a limit, we show

612 results with multiple values of Z_{max} in Supplemental Figure 3 and Supplemental Figure 5, while using a
613 default value of $Z_{max} = 6$ in other figures. In this manner, we have used the results for an infinite
614 system in which correlations are absent, but applied them to a system in which the number of units
615 could range from 10^3 to 10^6 to 10^{15} (as Z_{max} changes from 3 to 6 to 9) and the results be accurate for
616 999 networks in 1000 of that size. For further explanation see also the text in Section 3.1.
617

618 References

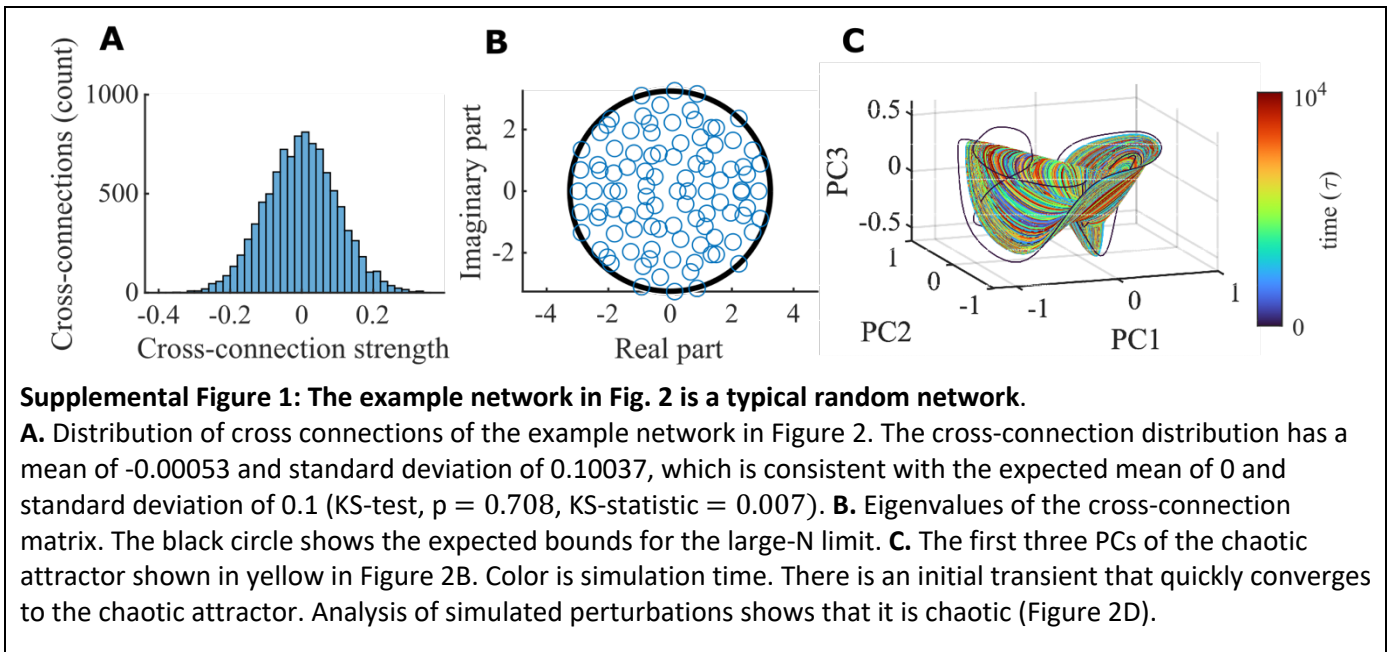
- 619
- 620 Abeles, M., Bergman, H., Gat, I., Meilijson, I., Seidemann, E., Tishby, N., & Vaadia, E. (1995).
621 Cortical activity flips among quasi-stationary states. *Proc Natl Acad Sci U S A*, *92*(19),
622 8616-8620.
- 623 Ahmadian, Y., Fumarola, F., & Miller, K. D. (2015). Properties of networks with partially
624 structured and partially random connectivity. *Phys Rev E Stat Nonlin Soft Matter Phys*,
625 *91*(1), 012820. <https://doi.org/10.1103/PhysRevE.91.012820>
- 626 Amit, D. J., Gutfreund, H., & Sompolinsky, H. (1985a). Spin-glass models of neural networks.
627 *Phys Rev A Gen Phys*, *32*(2), 1007-1018. <https://doi.org/10.1103/physreva.32.1007>
- 628 Amit, D. J., Gutfreund, H., & Sompolinsky, H. (1985b). Storing infinite numbers of patterns in a
629 spin-glass model of neural networks. *Phys. Rev. Lett.*, *55*, 1530-1531.
- 630 Anishchenko, A., & Treves, A. (2006). Autoassociative memory retrieval and spontaneous
631 activity bumps in small-world networks of integrate-and-fire neurons. *J Physiol Paris*,
632 *100*(4), 225-236. <https://doi.org/10.1016/j.jphysparis.2007.01.004>
- 633 Ballintyn, B., Shlaer, B., & Miller, P. (2019). Spatiotemporal discrimination in attractor networks
634 with short-term synaptic plasticity. *J Comput Neurosci*, *46*(3), 279-297.
635 <https://doi.org/10.1007/s10827-019-00717-5>
- 636 Battaglia, F. P., & Treves, A. (1998). Stable and rapid recurrent processing in realistic
637 autoassociative memories. *Neural Comput*, *10*(2), 431-450.
638 <http://www.ncbi.nlm.nih.gov/pubmed/9472489>
- 639 Benozzo, D., La Camera, G., & Genovesio, A. (2021). Slower prefrontal metastable dynamics
640 during deliberation predicts error trials in a distance discrimination task. *Cell Rep*, *35*(1),
641 108934. <https://doi.org/10.1016/j.celrep.2021.108934>
- 642 Boboeva, V., Pezzotta, A., & Clopath, C. (2021). Free recall scaling laws and short-term memory
643 effects in a latching attractor network. *Proc Natl Acad Sci U S A*, *118*(49).
644 <https://doi.org/10.1073/pnas.2026092118>
- 645 Bourjaily, M. A., & Miller, P. (2011). Excitatory, inhibitory, and structural plasticity produce
646 correlated connectivity in random networks trained to solve paired-stimulus tasks.
647 *Frontiers in Computational Neuroscience*, *5*, 37.
648 <https://doi.org/10.3389/fncom.2011.00037>
- 649 Brunel, N. (2003). Dynamics and plasticity of stimulus-selective persistent activity in cortical
650 network models. *Cereb Cortex*, *13*(11), 1151-1161.
651 [http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Cita
652 tion&list_uids=14576207](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list_uids=14576207)
- 653 Daelli, V., & Treves, A. (2010). Neural attractor dynamics in object recognition. *Exp Brain Res*,
654 *203*(2), 241-248. <https://doi.org/10.1007/s00221-010-2243-1>
- 655 David, H. A., & Nagaraja, H. N. (2003). *Order Statistics* (3rd ed.). John Wiley and Sons.
656 <https://doi.org/http://dx.doi.org/10.1002/0471722162>
- 657 Escola, S., Fontanini, A., Katz, D., & Paninski, L. (2011). Hidden Markov models for the stimulus-
658 response relationships of multistate neural systems. *Neural Comput*, *23*(5), 1071-1132.
659 https://doi.org/10.1162/NECO_a_00118

- 660 Folli, V., Leonetti, M., & Ruocco, G. (2016). On the Maximum Storage Capacity of the Hopfield
661 Model. *Front Comput Neurosci*, 10, 144. <https://doi.org/10.3389/fncom.2016.00144>
- 662 Fuster, J. M. (1973). Unit activity in prefrontal cortex during delayed-response performance:
663 neuronal correlates of transient memory. *Journal of Neurophysiology*, 36(1), 61-78.
664 <https://doi.org/10.1152/jn.1973.36.1.61>
- 665 Goldberg, J. A., Rokni, U., & Sompolinsky, H. (2004). Patterns of ongoing activity and the
666 functional architecture of the primary visual cortex. *Neuron*, 42(3), 489-500.
667 <http://www.ncbi.nlm.nih.gov/pubmed/15134644>
- 668 Golos, M., Jirsa, V., & Dauce, E. (2015). Multistability in Large Scale Models of Brain Activity.
669 *PLoS Comput Biol*, 11(12), e1004644. <https://doi.org/10.1371/journal.pcbi.1004644>
- 670 Hebb, D. O. (1949). *The organization of behavior; a neuropsychological theory*. Wiley.
- 671 Hopfield, J. J. (1982). Neural networks and physical systems with emergent collective
672 computational abilities. *Proceedings of the National Academy of Sciences*, 79(8), 2554-
673 2558. <https://doi.org/10.1073/pnas.79.8.2554>
- 674 Hopfield, J. J. (1984). Neurons with graded response have collective computational properties
675 like those of two-state neurons. *Proc. Natl. Acad. Sci. U.S.A.*, 81, 3088-3092.
- 676 Jones, L. M., Fontanini, A., Sadacca, B. F., Miller, P., & Katz, D. B. (2007). Natural stimuli evoke
677 dynamic sequences of states in sensory cortical ensembles. *Proc Natl Acad Sci U S A*,
678 104(47), 18772-18777. <https://doi.org/10.1073/pnas.0705546104>
- 679 Ksander, J., Katz, D. B., & Miller, P. (2021). A model of naturalistic decision making in preference
680 tests. *PLoS Comput Biol*, 17(9), e1009012. <https://doi.org/10.1371/journal.pcbi.1009012>
- 681 La Camera, G., Fontanini, A., & Mazzucato, L. (2019). Cortical computations via metastable
682 activity. *Curr Opin Neurobiol*, 58, 37-45. <https://doi.org/10.1016/j.conb.2019.06.007>
- 683 Lerner, I., Bentin, S., & Shriki, O. (2012). Spreading activation in an attractor network with
684 latching dynamics: automatic semantic priming revisited. *Cognitive science*, 36(8), 1339-
685 1382. <https://doi.org/10.1111/cogs.12007>
- 686 Lerner, I., Bentin, S., & Shriki, O. (2014). Integrating the automatic and the controlled: strategies
687 in semantic priming in an attractor network with latching dynamics. *Cognitive science*,
688 38(8), 1562-1603. <https://doi.org/10.1111/cogs.12133>
- 689 Lerner, I., & Shriki, O. (2014). Internally- and externally-driven network transitions as a basis for
690 automatic and strategic processes in semantic priming: theory and experimental
691 validation. *Frontiers in psychology*, 5, 314. <https://doi.org/10.3389/fpsyg.2014.00314>
- 692 Linkerhand, M., & Gros, C. (2013). Generating functionals for autonomous latching dynamics in
693 attractor relict networks. *Sci Rep*, 3, 2042. <https://doi.org/10.1038/srep02042>
- 694 Mazzucato, L., Fontanini, A., & La Camera, G. (2015). Dynamics of multistable states during
695 ongoing and evoked cortical activity. *J Neurosci*, 35(21), 8214-8231.
696 <https://doi.org/10.1523/JNEUROSCI.4819-14.2015>
- 697 Mazzucato, L., La Camera, G., & Fontanini, A. (2019). Expectation-induced modulation of
698 metastable activity underlies faster coding of sensory stimuli. *Nat Neurosci*, 22(5), 787-
699 796. <https://doi.org/10.1038/s41593-019-0364-9>
- 700 Miller, P. (2013). Stimulus number, duration and intensity encoding in randomly connected
701 attractor networks with synaptic depression. *Front Comput Neurosci*, 7, 59.
702 <https://doi.org/10.3389/fncom.2013.00059>

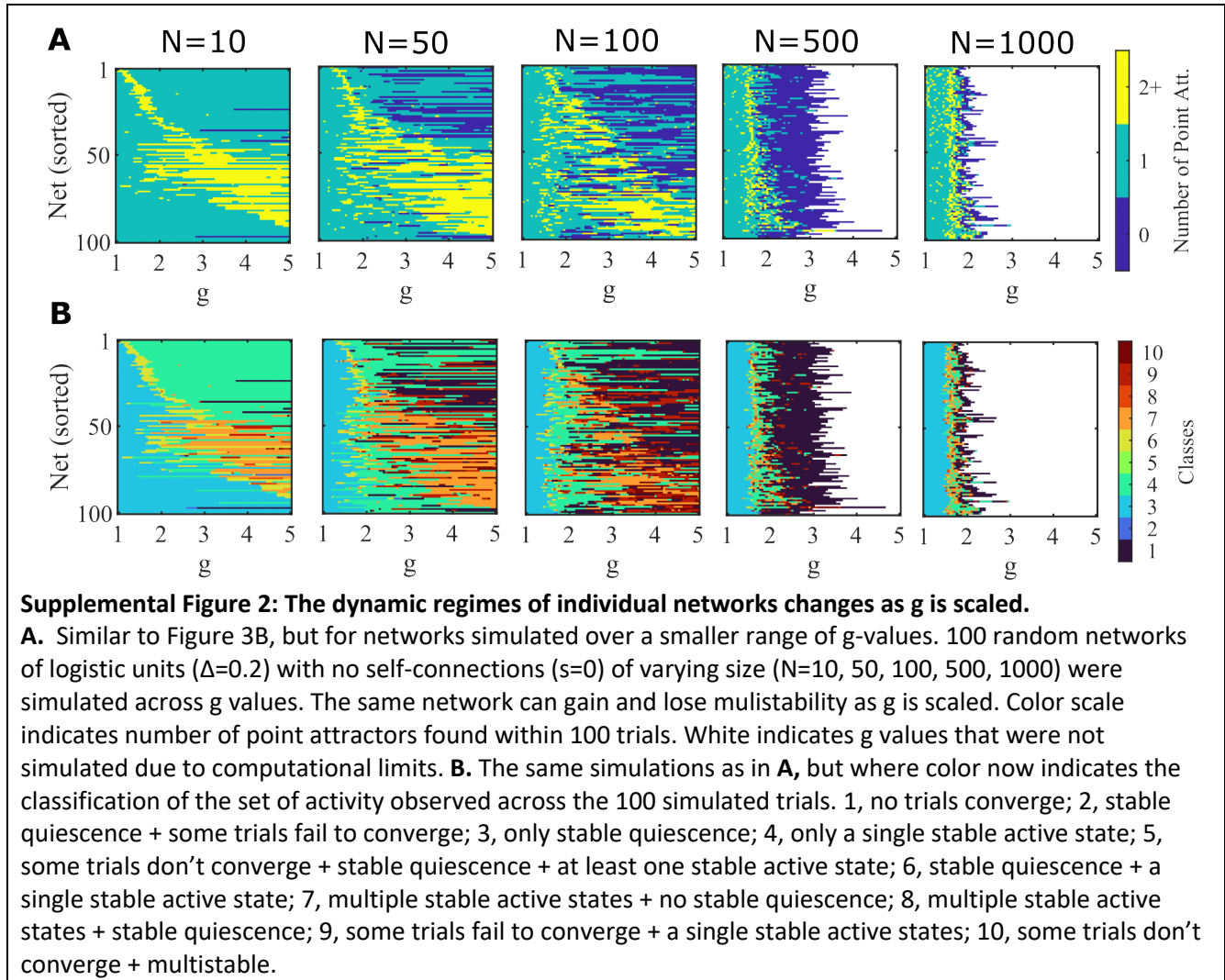
- 703 Miller, P. (2016). Itinerancy between attractor states in neural systems. *Curr Opin Neurobiol*, 40,
704 14-22. <https://doi.org/10.1016/j.conb.2016.05.005>
- 705 Miller, P., & Katz, D. B. (2010). Stochastic Transitions between Neural States in Taste Processing
706 and Decision-Making. *The Journal of Neuroscience*, 30(7), 2559-2570.
707 <https://doi.org/10.1523/jneurosci.3047-09.2010>
- 708 Miller, P., & Katz, D. B. (2011). Stochastic Transitions between States of Neural Activity. In M.
709 Ding & D. L. Glanzman (Eds.), *The Dynamic Brain: An Exploration of Neuronal Variability*
710 *and Its Functional Significance* (pp. 29-46). Oxford University Press.
- 711 Morcos, A. S., & Harvey, C. D. (2016). History-dependent variability in population dynamics
712 during evidence accumulation in cortex. *Nat Neurosci*, 19(12), 1672-1681.
713 <https://doi.org/10.1038/nn.4403>
- 714 Moreno-Bote, R., Rinzel, J., & Rubin, N. (2007). Noise-induced alternations in an attractor
715 network model of perceptual bistability. *J Neurophysiol*, 98(3), 1125-1139.
716 <https://doi.org/10.1152/jn.00116.2007>
- 717 Perin, R., Berger, T. K., & Markram, H. (2011). A synaptic organizing principle for cortical
718 neuronal groups. *Proc Natl Acad Sci U S A*, 108(13), 5419-5424.
719 <https://doi.org/10.1073/pnas.1016051108> [pii]
720 10.1073/pnas.1016051108
- 721 Ponce-Alvarez, A., Nacher, V., Luna, R., Riehle, A., & Romo, R. (2012). Dynamics of cortical
722 neuronal ensembles transit from decision making to storage for later report. *The Journal*
723 *of neuroscience : the official journal of the Society for Neuroscience*, 32(35), 11956-
724 11969. <https://doi.org/10.1523/JNEUROSCI.6176-11.2012>
- 725 Rabinovich, M., Volkovskii, A., Lecanda, P., Huerta, R., Abarbanel, H. D., & Laurent, G. (2001).
726 Dynamical encoding by networks of competing neuron groups: winnerless competition
727 [Research Support, Non-U.S. Gov't
728 Research Support, U.S. Gov't, Non-P.H.S.]. *Physical Review Letters*, 87(6), 068102.
729 <http://www.ncbi.nlm.nih.gov/pubmed/11497865>
- 730 Rabinovich, M. I., Varona, P., Tristan, I., & Afraimovich, V. S. (2014). Chunking dynamics:
731 heteroclinics in mind. *Front Comput Neurosci*, 8, 22.
732 <https://doi.org/10.3389/fncom.2014.00022>
- 733 Rainer, G., & Miller, E. K. (2000). Neural ensemble states in prefrontal cortex identified using a
734 hidden Markov model with a modified EM algorithm. *Neurocomputing*, 32, 961-966.
735 [https://doi.org/Doi 10.1016/S0925-2312\(00\)00266-6](https://doi.org/Doi%2010.1016/S0925-2312(00)00266-6)
- 736 Rajan, K., & Abbott, L. F. (2006). Eigenvalue spectra of random matrices for neural networks.
737 *Phys Rev Lett*, 97(18), 188104. <http://www.ncbi.nlm.nih.gov/pubmed/17155583>
- 738 Recanatesis, S., Pereira, U., Murakami, M., Mainen, Z., & Mazzucato, L. (2022). Metastable
739 attractors explain the variable timing of stable behavioral action sequences. *Neuron*.
- 740 Russo, E., & Treves, A. (2012). Cortical free-association dynamics: distinct phases of a latching
741 network. *Phys Rev E Stat Nonlin Soft Matter Phys*, 85(5 Pt 1), 051920.
742 <https://doi.org/10.1103/PhysRevE.85.051920>
- 743 Sadacca, B. F., Mukherjee, N., Vladusich, T., Li, J. X., Katz, D. B., & Miller, P. (2016). The
744 Behavioral Relevance of Cortical Neural Ensemble Responses Emerges Suddenly. *The*
745 *Journal of Neuroscience*, 36(3), 655-669. [https://doi.org/10.1523/jneurosci.2265-](https://doi.org/10.1523/jneurosci.2265-15.2016)
746 [15.2016](https://doi.org/10.1523/jneurosci.2265-15.2016)

- 747 Seidemann, E., Meilijson, I., Abeles, M., Bergman, H., & Vaadia, E. (1996). Simultaneously
748 recorded single units in the frontal cortex go through sequences of discrete and stable
749 states in monkeys performing a delayed localization task. *J Neurosci*, *16*(2), 752-768.
750 <http://www.ncbi.nlm.nih.gov/pubmed/8551358>
- 751 Sompolinsky, H., & Crisanti, A. (2018). Path integral approach to random neural networks.
752 *Physical Review E*, *98*, 062120.
- 753 Sompolinsky, H., Crisanti, A., & Sommers, H. J. (1988). Chaos in random neural networks. *Phys*
754 *Rev Lett*, *61*(3), 259-262. <https://doi.org/10.1103/PhysRevLett.61.259>
- 755 Song, S., Sjöström, P. J., Reigl, M., Nelson, S., & Chklovskii, D. B. (2005). Highly nonrandom
756 features of synaptic connectivity in local cortical circuits. *PLoS Biol*, *3*(3), e68.
757 <https://doi.org/10.1371/journal.pbio.0030068> [pii]
- 758 10.1371/journal.pbio.0030068
- 759 Song, S., Yao, H., & Treves, A. (2014). A modular latching chain. *Cogn Neurodyn*, *8*(1), 37-46.
760 <https://doi.org/10.1007/s11571-013-9261-1>
- 761 Stepanyants, A., & Chklovskii, D. B. (2005). Neurogeometry and potential synaptic connectivity.
762 *Trends Neurosci*, *28*, 387-394.
- 763 Stern, M., Sompolinsky, H., & Abbott, L. F. (2014). Dynamics of random neural networks with
764 bistable units. *Physical review. E, Statistical, nonlinear, and soft matter physics*, *90*(6),
765 062710-062710. <https://doi.org/10.1103/PhysRevE.90.062710>
- 766 Strogatz, S. H. (2015). *Nonlinear Dynamics and Chaos* (2nd ed.). Westview Press.
- 767 Taylor, J. D., Chauhan, A. S., Taylor, J. T., Shilnikov, A. L., & Nogaret, A. (2022). Noise-activated
768 barrier crossing in multiattractor dissipative neural networks. *Phys Rev E*, *105*(6-1),
769 064203. <https://doi.org/10.1103/PhysRevE.105.064203>
- 770 Touboul, J. D., & Ermentrout, G. B. (2011). Finite-size and correlation-induced effects in mean-
771 field dynamics. *J Comput Neurosci*, *31*(3), 453-484. [https://doi.org/10.1007/s10827-011-](https://doi.org/10.1007/s10827-011-0320-5)
772 [0320-5](https://doi.org/10.1007/s10827-011-0320-5)
- 773 Treves, A. (1990). Graded-response neurons and information encodings in autoassociative
774 memories. *Phys Rev A*, *42*(4), 2418-2430.
775 <http://www.ncbi.nlm.nih.gov/pubmed/9904294>
- 776 Treves, A. (2005). Frontal latching networks: a possible neural basis for infinite recursion. *Cogn*
777 *Neuropsychol*, *22*(3), 276-291. <https://doi.org/10.1080/02643290442000329>
- 778 Wills, T. J., Lever, C., Cacucci, F., Burgess, N., & O'Keefe, J. (2005). Attractor Dynamics in the
779 Hippocampal Representation of the Local Environment. *Science*, *308*(5723), 873-876.
780 [https://doi.org/doi:10.1126/science.1108905](https://doi.org/10.1126/science.1108905)
- 781 Wilson, H., & Cowan, J. (1973). A Mathematical Theory of the Functional Dynamics of Cortical
782 and Thalamic Nervous Tissue. *Kybernetik*, *13*, 55-80.
783 <https://doi.org/10.1007/BF00288786>
- 784 Zurada, J. M., Cloete, I., & van der Poel, E. (1996). Generalized Hopfield networks for associative
785 memories with multi-valued stable states. *Neurocomputing*, *13*, 135-149.
786
- 787
- 788

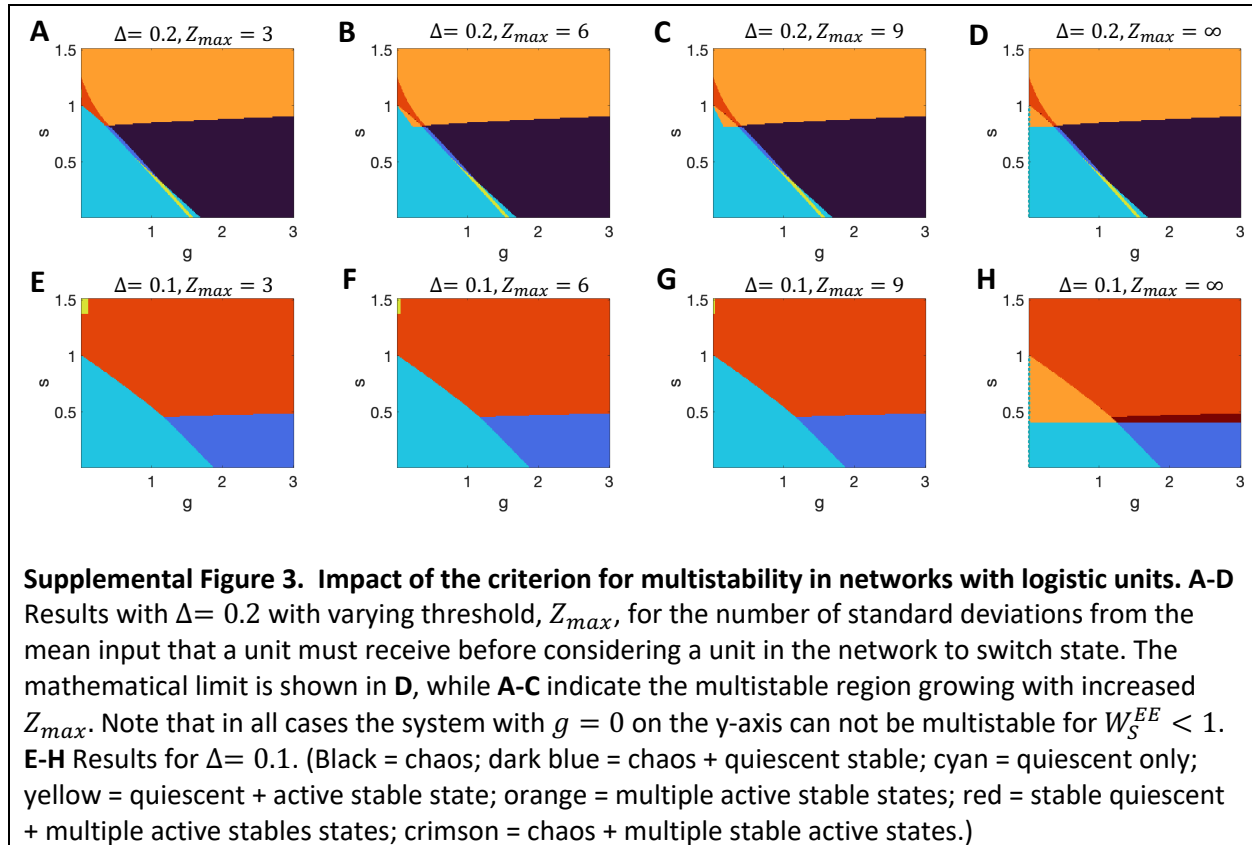
789 Supplemental figures:



790

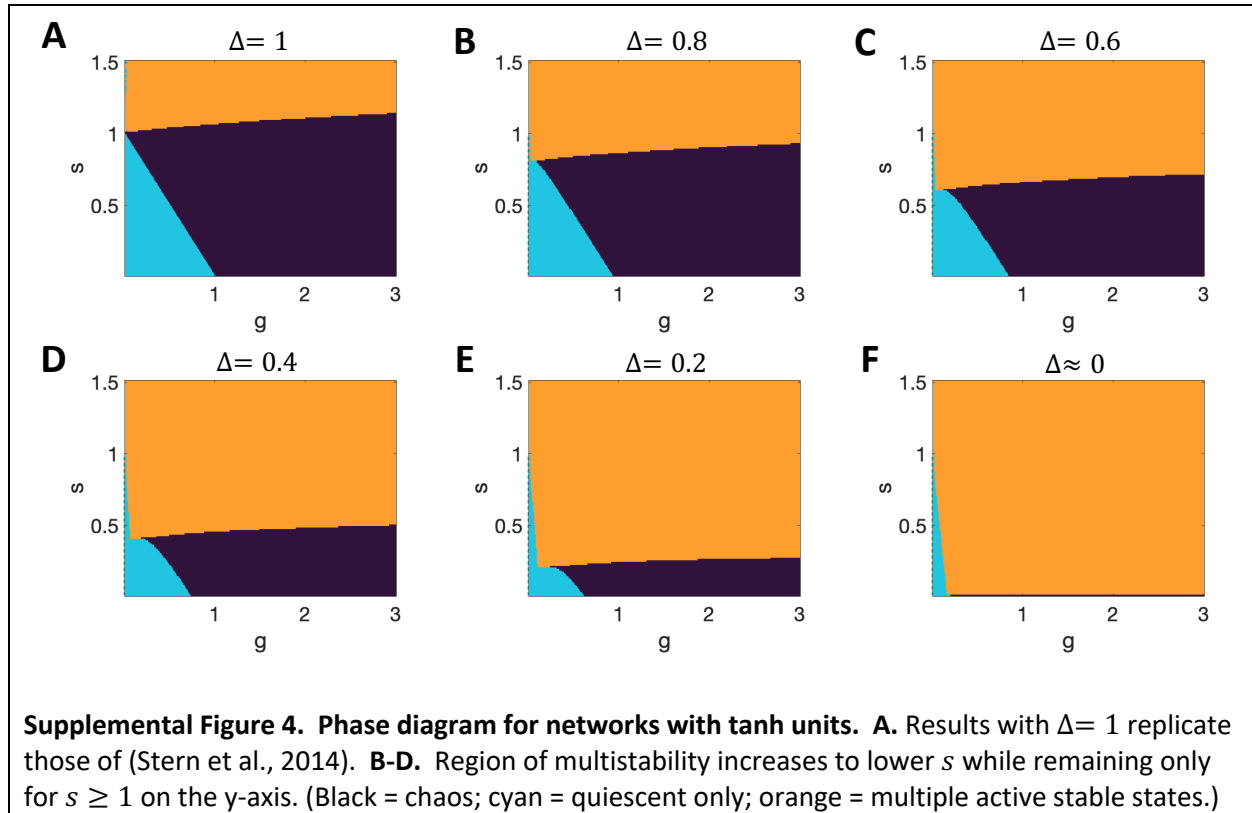


791
792

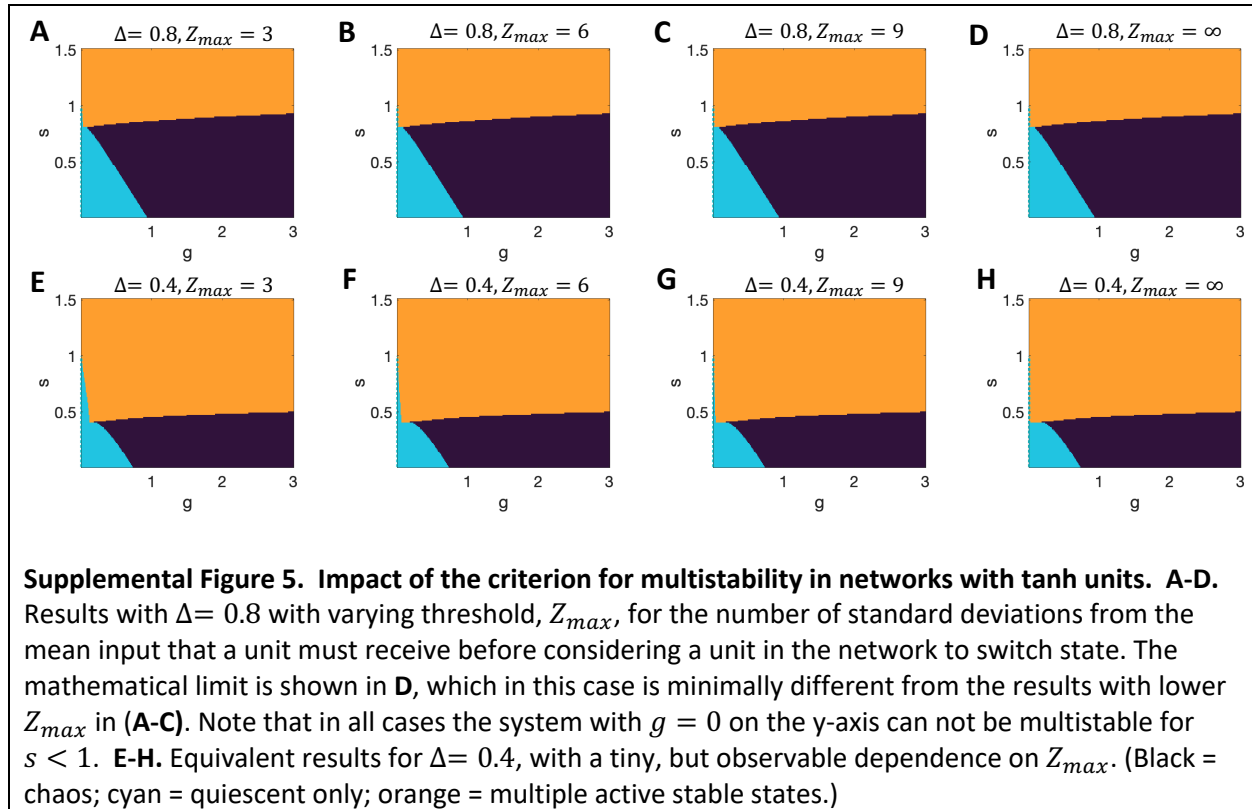


793

794



795
796



797