

Supporting Information

Tkačik et al. 10.1073/pnas.1514188112

Thermodynamics and Probability Distributions

The fundamental variables of thermodynamics are energy, temperature, and entropy. For the states taken on by a network of neurons, energy and temperature are meaningless, so it is difficult to see how we can construct a thermodynamics for these systems. However, in statistical mechanics, all thermodynamic quantities are derivable from the Boltzmann distribution, the probability that the system will be found in any particular state. Thus, all thermodynamic statements are equivalent to statements about this underlying probability distribution, and, in this sense, we should be able to construct thermodynamics for a much broader range of probability distributions that describe a large number of variables.

The idea that probability distributions over N variables can have an associated thermodynamics in the $N \rightarrow \infty$ limit is powerful but perhaps not so widely used. This connection is well-studied by mathematical physicists (1) and has been a guide to the analysis of experiments on dynamical systems (2, 3). We have used these ideas to construct a thermodynamics of natural images (4) and have emphasized the connection of thermodynamic criticality to Zipf's law (5). Here we give a somewhat pedagogical discussion, in the hope of making the results accessible to a broader audience.

The Boltzmann Distribution. We start by recalling that, for a system in thermal equilibrium at temperature T , the probability of finding the system in state s is given by

$$P_s = \frac{1}{Z} e^{-E_s/k_B T}, \quad [\text{S1}]$$

where E_s is the energy of the state, and Boltzmann's constant k_B converts between conventional units of temperature and energy. The partition function Z serves to normalize the distribution, which requires

$$Z = \sum_s e^{-E_s/k_B T}, \quad [\text{S2}]$$

but, in fact, this normalization constant encodes many physical quantities. The logarithm of the partition function is proportional to the free energy of the system, the derivative of the free energy with respect to the volume occupied by the system is the pressure, the derivative with respect to the strength of an applied magnetic field is the magnetization, and so on.

The state of a system is defined by the joint configuration of all its parts. Thus, in a classical gas or liquid, s is defined by the positions and velocities of all of the constituent atoms. Different gases or liquids differ not because these variables are different but because the energy E_s is a different function of these N underlying variables. However, thermodynamics doesn't make reference to all these details. Which aspects of the underlying microscopic rules actually matter for predicting the free energy and its derivatives?

We can write the sum over all states as a sum first over states that have the same energy and then as a sum over energies. We do this by introducing an integral over a delta function into the sum,

$$Z = \sum_s e^{-E_s/k_B T} = \sum_s \left[\int dE \delta(E - E_s) \right] e^{-E_s/k_B T} \quad [\text{S3}]$$

$$= \int dE \sum_s \delta(E - E_s) e^{-E_s/k_B T} \quad [\text{S4}]$$

$$= \int dE e^{-E/k_B T} \left[\sum_s \delta(E - E_s) \right]. \quad [\text{S5}]$$

We see that the way in which the energy depends on each state appears only in the brackets, a function $n(E)$ that counts how many states have a particular energy,

$$n(E) = \sum_s \delta(E - E_s). \quad [\text{S6}]$$

Looking ahead to the analysis of real data, it will be convenient to rearrange Eq. S5 slightly. Instead of counting the number of states that have energy E , we can count the number of states with energy less than E ,

$$\mathcal{N}(E) = \sum_s \Theta(E - E_s), \quad [\text{S7}]$$

where the step function is defined by

$$\Theta(x > 0) = 1 \quad [\text{S8}]$$

$$\Theta(x < 0) = 0. \quad [\text{S9}]$$

However, the step function is the integral of the delta function, which means that we can integrate by parts in Eq. S5 to give

$$Z = \frac{1}{k_B T} \int dE e^{-E/k_B T} \mathcal{N}(E). \quad [\text{S10}]$$

If we think about N variables, each of which can take on only two states, the total number of states is 2^N . More generally, we expect that the number of possible states in a system with N variables is exponentially large, so it is natural to think not about the number of states $\mathcal{N}(E)$ but about its logarithm,

$$S(E) = \ln \mathcal{N}(E), \quad [\text{S11}]$$

which is called the entropy.

As a technical aside, we can either define the entropy in terms of the density of states with energy close to E , what we have called $n(E)$, or use the number of states with energy less than E , what we have called $\mathcal{N}(E)$. When the number of degrees of freedom is small, these are both badly behaved functions— $n(E)$ is singular, and $\mathcal{N}(E)$ has visible steps. However, as N becomes large, both functions become smooth, and we can do all of the usual operations of differentiation or integration by parts without worries. Importantly, when it comes time to analyze experimental data, using $\mathcal{N}(E)$ allows us to avoid making bins along the E axis.

Substituting from Eq. S11 into Eq. S10, the partition function can be written as an integral determined only by the function $S(E)$, entropy vs. energy,

$$Z = \frac{1}{k_B T} \int dE \exp\left[-\frac{E}{k_B T} + S(E)\right]. \quad [\text{S12}]$$

One of the key ideas in thermodynamics is that certain variables are “extensive,” that is, proportional to the number of particles or variables in the system, whereas other variables are “intensive,” independent of the system size. Temperature is an intensive variable, whereas energy and entropy are extensive variables. It is then natural to think about the energy per particle, $\epsilon = E/N$, and the entropy per particle, $S(E)/N = s(\epsilon)$. In the limit of large N , we expect $s(\epsilon)$ to become a smooth function. Substituting into Eq. S12, the partition function can be written as

$$Z = \frac{N}{k_B T} \int d\epsilon e^{-Nf(\epsilon)/k_B T} \quad [\text{S13}]$$

$$f(\epsilon) = \epsilon - k_B T s(\epsilon). \quad [\text{S14}]$$

We note that $f(\epsilon)$ is the difference between energy and entropy, scaled by the temperature, and is called the free energy.

Whenever we have an integral of the form in Eq. S13, at large N , we expect that it will be dominated by values of ϵ close to the minimum of $f(\epsilon)$. This minimum ϵ_* is the solution to the equation

$$\frac{df(\epsilon)}{d\epsilon} = 0 \Rightarrow \frac{1}{k_B T} = \frac{ds(\epsilon)}{d\epsilon}, \quad [\text{S15}]$$

which we can also think of as defining the temperature. Notice that T being positive requires that the system have $ds(\epsilon)/d\epsilon > 0$, which means there are more states with higher energies.

If we expand $f(\epsilon)$ in the neighborhood of ϵ_* , we have

$$f(\epsilon) = f(\epsilon_*) - \frac{k_B T}{2} \frac{d^2 s(\epsilon)}{d\epsilon^2} \Big|_{\epsilon_*} (\epsilon - \epsilon_*)^2 + \dots, \quad [\text{S16}]$$

which gives

$$Z \approx \frac{N}{k_B T} e^{-Nf(\epsilon_*)/k_B T} \int d\epsilon e^{-A(\epsilon - \epsilon_*)^2} \quad [\text{S17}]$$

$$A = \frac{N}{2} \left[\frac{d^2 s(\epsilon)}{d\epsilon^2} \Big|_{\epsilon_*} \right]. \quad [\text{S18}]$$

This looks as if the energy per particle is drawn from an approximately Gaussian distribution, with mean ϵ_* and variance

$$\langle (\delta\epsilon)^2 \rangle = \frac{1}{N} \left[\frac{d^2 s(\epsilon)}{d\epsilon^2} \Big|_{\epsilon_*} \right]^{-1}, \quad [\text{S19}]$$

and, indeed, this can be shown more directly from the Boltzmann distribution.

With the interpretation of ϵ_* as the mean energy per particle, we can use Eq. S15 to calculate how this energy changes when we change the temperature, and we find

$$\frac{d\epsilon_*}{dT} = \frac{1}{k_B T^2} \left[-\frac{d^2 s(\epsilon)}{d\epsilon^2} \Big|_{\epsilon_*} \right]^{-1}. \quad [\text{S20}]$$

The change in energy with temperature is called the heat capacity C , and, when we normalize per particle, it is referred to as the

specific heat. Combining Eqs. S19 and S20, we see that the specific heat C/N is connected to the variance in energies,

$$\langle (\delta\epsilon)^2 \rangle = k_B T^2 \frac{C}{N}. \quad [\text{21}]$$

This relationship also can be proven without resorting to the approximation in Eq. S16.

Our discussion thus far assumes that the second derivative of the entropy with respect to the energy is not zero. If we take all our results at face value, then, when $d^2 s/d\epsilon^2 \rightarrow 0$, the specific heat will become infinite (Eq. S20), as will the variance of the energy per particle (Eq. S19). This is a critical point.

There is much more to be said about the analysis of critical points using the entropy vs. energy. However, our concern here is how these ideas connect to systems that are not in thermal equilibrium, so that temperature and energy are not relevant concepts. What we would like to show is that many of the thermodynamic quantities nonetheless serve to characterize the behavior of any probability distribution for a very large number of variables.

Distributions, More Generally. Rather than trying to compute the partition function, we can ask, for any distribution, how the normalization condition is satisfied. We still imagine that there are states s , built of N different variables, as with the patterns of spiking and silence in a network of neurons. Each state s has a probability P_s , and we must have

$$1 = \sum_s P_s. \quad [\text{S22}]$$

We can now follow the same strategy that we used above for the partition function: We do the sum first by summing over all of the states that have the same value of the (log) probability, and then we sum over this value. We start by defining

$$E_s = -\ln P_s, \quad [\text{S23}]$$

as in Eq. 1. Then we have

$$\sum_s P_s = \sum_s \int dE \delta(E - E_s) P_s. \quad [\text{S24}]$$

However, because $P_s = e^{-E_s}$, we can rewrite this as

$$\sum_s P_s = \int dE e^{-E} \sum_s \delta(E - E_s). \quad [\text{S25}]$$

Integrating by parts, we obtain

$$\sum_s P_s = \int dE e^{-E} \mathcal{N}(E), \quad [\text{S26}]$$

where $\mathcal{N}(E)$ is a cumulative density of states, as in Eq. S7,

$$\mathcal{N}(E) = \sum_s \Theta(E - E_s). \quad [\text{S27}]$$

Again, this is a number of states, so the logarithm of this number is an entropy, exactly as in Eq. S11. Thus, the statement that the probability distribution is normalized becomes

$$\sum_s P_s = \int dE \exp[-E + S(E)]. \quad [\text{S28}]$$

If we have a system in which the state s is built out of N variables, then we expect that, for large N , both log probabilities (E) and entropies (S) are proportional to N . A standard example is in information theory, where s could label a message built out of N symbols, and the proportionality $E \propto N$ is central to proofs of the classic coding theorems (6). In the case of interest to us here, we can look at the states taken on by groups of N neurons, and we can vary N over some range. The function $\mathcal{N}(E)$, and hence the entropy $S(E)$, is a property of a single system with a particular value of N , and, to remind us of this fact, we can write $S_N(E)$. What happens as N become large is an experimental question. However, in many of the examples that we understand—from statistical physics, from information theory, and indeed from more general examples in probability theory—there is a well-defined limiting behavior at large N , which means that there is a function

$$s(\epsilon) = \lim_{N \rightarrow \infty} \frac{1}{N} S_N(E = N\epsilon). \quad [\text{S29}]$$

If this limit exists, then the normalization condition on the probability distribution in Eq. S28 becomes

$$\sum_s P_s \rightarrow NP_0 \int d\epsilon e^{-Nf(\epsilon)}, \quad [\text{S30}]$$

$$f(\epsilon) = \epsilon - s(\epsilon). \quad [\text{S31}]$$

Now we can see the correspondence with the description of an equilibrium thermodynamic system, which leads up to the expressions for the partition function in Eqs. S13 and S14:

- i) We can assign an energy to every state of the system, which is just the negative log probability. The effective temperature of the system is $k_B T = 1$.
- ii) We can count the number of states below a given energy, and the log of this number is an entropy.
- iii) If there are N elements (e.g., neurons) in our system, it is natural to ask about the entropy per element as a function of the energy per element. If this function has a smooth limit as N becomes large, $s(\epsilon)$, then we can define a thermodynamics for the system.
- iv) When we sum over states, the sum is dominated by states that minimize the free energy, $f(\epsilon) = \epsilon - s(\epsilon)$, just as in ordinary thermodynamics, provided that the curvature of the free energy at this minimum is nonzero.
- v) The dominance of states near the minimum of the free energy enforces the notion of “typicality” (6), so that at large N most of the states we actually see have essentially the same value of log probability.
- vi) If the curvature at the minimum of the free energy vanishes, then the usual ideas of typicality break down, and we will see large fluctuations in the log probability of states, even if we normalize this log probability by N .
- vii) The large variance in log probability is mathematically equivalent to a diverging specific heat in the thermodynamic case. This is a signature of a critical point.

More About Criticality. Before leaving this discussion, we note that there are other signatures of criticality, and even different notions of criticality. In equilibrium systems with interactions that extend only over short distances, correlations typically extend over some longer but finite distance ξ ; at the critical point, this correlation length diverges, so that there is no characteristic length scale—all scales between the size of the constituent particles and the size of the system as a whole are relevant (7). Not only does the specific heat diverge at the critical point, but so does the susceptibility to external fields. All of these diverging quantities

have a power law dependence on the difference between the actual temperature and the critical temperature, and the exponents of these power laws are quantitatively universal: Many different systems, with different microscopic constituents, exhibit precisely the same exponents, and, in a certain precise sense, these exponents give a complete description of the system in the neighborhood of the critical point (8, 9). In the study of complex, nonequilibrium systems, scale invariance and power law behaviors often are taken as signs of criticality, but seldom is it possible to exhibit these behaviors over the wide range of scales that are the standard in studies of equilibrium critical phenomena, so one must be cautious.

In almost all equilibrium systems, the approach to criticality also is associated with the emergence of long time scales in the dynamics; as with the divergence of the correlation length ξ , the divergence of the correlation time in the dynamics means that there is a form of temporal scale invariance at criticality. Deterministic dynamical systems also exhibit critical phenomena, often called bifurcations, where the system’s behavior changes qualitatively in response to an infinitesimal change in parameters (10). These phenomena are easiest to understand when the number of degrees of freedom N is small, but then the sharp bifurcations are rounded if there is noise in the system; the example of equilibrium statistical mechanics shows how noisy dynamical systems can recover sharp transitions in the limit of large N . In general, it is not clear how dynamical and statistical notions of criticality are related to one another in systems with many degrees of freedom.

Experimental Methods

Much of the analysis in this paper is based on the same data set as in ref. 11. For completeness, we review our experimental methods here. Retinae were isolated from the eye in darkness, and the retina was pressed against a custom-fabricated array of 252 electrodes. The retina was superfused with oxygenated Ringer’s medium at room temperature. Electrode voltage signals were acquired and digitized at 10 kHz by a 252-channel preamplifier (MultiChannel Systems). The sorting of these signals into action potentials from individual neurons was done offline using the methods of ref. 12.

The repeated natural movie was a movie of a fish tank captured at 30 Hz with a standard camera; it lasted 20 s and was repeated 297 times. As noted in the main text, this experiment allowed us to resolve 160 neurons across the recording array. The random checkerboard consisted of square pixels, 69 μm on a side, each chosen independently black or white 30 times per second, creating a 30-s random movie that was repeated 69 times; this experiment yielded 120 stable, resolved cells. For the spatially uniform flicker, the luminance of the entire screen was chosen randomly from a Gaussian distribution 60 times per second, creating a 10-s-long random sequence that was repeated 98 times; we separated the signals from 111 neurons.

Effects of Bin Size

We follow earlier work and define the states of the neural network in discrete time bins (13). That is, we slice the time axis in bins of duration $\Delta\tau$, and define $\sigma_i = 1$ at time t if neuron i spikes in the window $[t, t + \Delta\tau)$, and $\sigma_i = 0$ otherwise. We choose $\Delta\tau = 20$ ms because this captures the structure of the correlation functions, but it should be admitted that there is some arbitrariness here. If we make bins too large, surely we are grouping together distinct responses of the network, whereas, if we make the bins too small, then meaningful correlations are spread over multiple bins, and we need to analyze the distribution of state sequences rather than instantaneous states (14).

One might hope, however, that there is a range of bin sizes over which the basic structure of the distribution $P(\{\sigma_i\})$ is constant. We test this in Fig. S1, which should be compared with Fig. 3. Fig. S1 shows the entropy vs. energy, computed directly from the data, with bin widths of $\Delta\tau = 10$ ms and 40 ms, whereas we use $\Delta\tau = 20$ ms in the main text. Although details vary a bit, in all

cases, we see the approach to $s(\epsilon) = \epsilon$ as N becomes large. Although much remains to be understood about the dynamics of the states in this network, Fig. S1 demonstrates that our main results do not depend sensitively on the choice of $\Delta\tau$.

Analysis of Maximum Entropy Models

We can take our maximum entropy model seriously as a statistical mechanics problem and use Monte Carlo simulation to generate samples of the states $\{\sigma_i\}$ drawn from our model distribution. Heat capacity curves were estimated by running a Metropolis Monte Carlo sampler independently at every T . Because the model assigns an energy $E = E(\{\sigma_i\})$ to each state, we can compute the mean and variance of E from a single long Monte Carlo run, and thus estimate the heat capacity through the thermodynamic identity in Eq. S21. Samples of the energy were collected at every sweep (roughly N spin flips); 2×10^6 sweeps were performed for every T .

To estimate the function $n(E)$ in the maximum entropy models, including the α -ensembles in Fig. 5, we used Wang–Landau sampling (15). In detail, the complete energy range was divided into 2×10^4 equidistant energy bins (6×10^3 for the α -ensembles), the histogram flatness criterion was 0.9, and the final multiplicative update was $1 + 10^{-5}$. These measurements, as well as the specific heat curves, can both be used to give an estimate of the entropy of the distribution, and these agree to within better than 1% (11), providing a check on our sampling procedures. For more on these matters, see the methods section, “Computing the entropy . . .,” of ref. 11.

In addition to the models described in the main text, we have also considered models that do not include the term $V(K)$ in Eq. 4; these are maximum entropy models that match exactly the mean spike probabilities of individual neurons, and the pairwise correlations, but not the probability of K neurons spiking simultaneously. As explained in ref. 11, these simpler, purely pairwise models provide noticeably less accurate descriptions of the network activity. [Note that the parameters $\{h_i; J_{ij}\}$ in the two models are not the same, but must be found, independently, to match the relevant expectation values (11).] Importantly, although both models capture all of the pairwise correlations among neurons, the peak of the specific heat is much stronger and more clearly N -dependent in the more accurate model, as shown in Fig. S2.

More About Alternatives

In this section, we expand on alternative interpretations of the data, arguing that the signatures of criticality are unlikely to be explained away as spurious consequences of less interesting models.

Impact of Limited Data. We have tested in detail the reliability with which maximum entropy models can be inferred from the available data. As explained in ref. 11, we can learn these models from 90% of the data and then compare the quality of the model against both the training set and the held-out test set. Even with $N = 120$ neurons, the model predicts that the log likelihood of the test data is the same as that of the training data, within error bars, and these errors are less than 1% (figure 4 of ref. 11). Still, one might worry that small errors associated with the finiteness of the data set could have a disproportionate impact on the putative signatures of criticality. To test for this, we have learned models for $N = 100$ neurons from fractions of the data ranging down to just 10%; results for the heat capacity vs. temperature (as in Fig. 4) are shown in Fig. S3. We see that the sharp peak in $C(T)$ is essentially independent of the sample size across this wide range, and that the variations in $C(T)$ across different small fractions of the data are only a few percent. Thus, this behavior is not a result of overfitting, nor is it linked in any way to the size of our data set.

It is important that, in Fig. S3, we are always looking at the same 100 neurons; otherwise, variability across subsets could be confused with sampling errors. When we change the size of the data, we are choosing, at random, some fraction of the experiment, and, for each fraction, we examine 10 such random choices. For each

choice, we make a completely independent reconstruction of the maximum entropy model, which means that variability includes not just the effects of finite data but also any errors in parameter estimation or in the Monte Carlo estimate of the specific heat. Evidently, all of these errors are quite small.

Are Correlations Inherited from the Visual Stimulus? As discussed in the main text, one possible interpretation of our observations is that correlations among neurons simply reflect correlations in the visual stimulus. In this case, any interesting features in the joint distribution of activity among many neurons would be entirely traceable to the structure of the sensory inputs.

The idea that correlations among neurons should be decomposed into contributions from their inputs and contributions intrinsic to the circuit is very old (16), dating back to a time when it was hoped that measurement of correlations would allow a direct inference of connectivity in the circuit. Before discussing the origin of correlations, it is important to emphasize that the distinction between “stimulus-induced” and “intrinsic” correlations is not a distinction that the brain can make. Experimentally, we make this distinction by providing exact repetitions of the stimulus, but this never happens in the natural world. The only knowledge that the brain has of its visual inputs is the set of signals provided by the population of ganglion cells itself, so there is no way to search for correlations with some other reference signal. We also note that, following decades of experiments on correlations among RGCs (17, 18), there is now direct evidence that triggering spikes in one ganglion cell changes the response of other ganglion cells to sensory signals,* so that these cells certainly are not responding independently to their visual inputs.

Although the dissection of the correlations is irrelevant for brain function, it is interesting to ask, mechanistically, how these correlations arise. If they arise solely from the visual inputs, then changing the statistical structure of these inputs should produce a dramatic effect. We have replaced the natural movies with randomly flickering checkerboards (an approximation to spatiotemporal white noise) and spatially uniform but temporally random flicker. In each case, we have constructed maximum entropy models (Eqs. 3 and 4) and searched for a peak in the specific heat vs. temperature, as in Fig. 4; results are shown in Fig. S4.

Although there are quantitative differences among the responses to the different stimulus ensembles, we see that there are signatures of criticality in each case. As with the natural movies, there is a peak in the specific heat, the height of the peak grows with the number of neurons, and the location of the peak moves toward $T = 1$ at larger N . It thus seems unlikely that these signatures of criticality in the specific heat are merely a reflection of input statistics. Indeed, we should remember that the decomposition of correlations into intrinsic and stimulus-induced is incomplete, because the retina adapts to the distribution of its inputs, on many time scales. It would appear that some combination of anatomical connectivity and adaptation poises the population of RGCs near a peak in the specific heat. This points toward future experiments that should probe more directly the invariance of thermodynamic behavior across adaptation states.†

It seems worth emphasizing that, even if correlations are largely inherited from the visual stimulus, this transformation from input to output correlations is nontrivial. The conventional model for the input–output relations of the neurons is the “linear–nonlinear” model, in which the probability of spiking is determined by an instantaneous nonlinear function of a linearly filtered version of the stimulus. In the retina, with the stimulus given by the light

*Asari H, Meister M, Computational and Systems Neuroscience, February 23–26, 2012, Salt Lake City, UT.

†Ioffe M, Tkačik G, Bialek W, Berry MJ, II, Computational and Systems Neuroscience, February 27 to March 2, 2014, Salt Lake City, UT.

intensity as a function of space and time, $I(\vec{x}, t)$, the probability of spiking for one cell in one time bin is then

$$p(t) = p_0 g \left[\int d^2x \int d\tau F(\vec{x}, \tau) I(\vec{x}, t - \tau) \right], \quad [\text{S32}]$$

where p_0 sets the maximum response, $g[\cdot]$ is a nonlinear function that we can normalize to range between 0 and 1, and $F(\vec{x}, \tau)$ is the linear spatiotemporal receptive field of the cell. It is a theorem that, if this model is an accurate description of the neural response, then the receptive field can be determined by correlating the spiking output with the spatiotemporal variations in the input, provided that the inputs are chosen from a Gaussian ensemble. In particular, if the inputs are white noise (down to the spatial and temporal resolution used), as in the random checkerboard experiments, then

$$F(\vec{x}, \tau) \propto \langle I(\vec{x}, t - \tau) \delta(t - t_{\text{spike}}) \rangle, \quad [\text{S33}]$$

where t_{spike} is the time of a spike and the average $\langle \dots \rangle$ is computed across a long sample of the random checkerboard movie. As described in the supporting information of ref. 19, we have constructed these receptive fields for every cell in the population, and then mapped the nonlinearities $g[\cdot]$ independently for each cell. If we then generate spikes at the output of this population, and compute their pairwise correlation coefficients, we obtain the results shown in Fig. S5.

There is a widely held intuition that correlations among neural responses in the retina should be understood as being shaped largely by the overlap of receptive fields, but Fig. S5 suggests that the situation is more complex. The linear–nonlinear model predicts correlations based on receptive field structure, but these predictions are strongly at variance with what we see in the data. The distribution of correlations in the model is narrower than in the data, failing to access the tail of strong positive correlations and cutting off at only modest negative correlations. Taking each pair of neurons individually, we see that the predicted and observed correlations are almost unrelated to one another.

Zipf's Law, Superposition, and Related Matters. Rather than counting states that have a particular value of the (log) probability, we can simply put the states in order of their probability, highest probability states first. The resulting plot of probability vs. rank is sometimes called the “Zipf plot,” with reference to the corresponding analysis of words in written language (20). As we have emphasized elsewhere (4, 5, 21), the Zipf plot is essentially the plot of entropy vs. energy, turned on its side. Concretely, if the state with rank r has probability p_r , then we have r states with probability $p \geq p_r$, or an effective energy $E \leq -\ln p_r$. However, the number of states with energy less than E is what we have called the cumulative density of states, $\mathcal{N}(E)$. Thus, we have

$$\mathcal{N}(E)|_{E=-\ln p_r} = r, \quad [\text{S34}]$$

or, for the entropy,

$$S(E)|_{E=-\ln p_r} = \ln r. \quad [\text{S35}]$$

What Zipf observed about words is that $p_r \approx A/r$, up to some maximum r . If we take this as an exact statement (“Zipf's law”), then $r = A/p_r$, and hence Eq. S35 becomes

$$S(E) = E + \ln A. \quad [\text{S36}]$$

Thus, Zipf's law is equivalent to a linear relation between entropy and energy, with slope one. Because this means that the second

derivative of the entropy with respect to the energy vanishes, Zipf's law seems to imply criticality, in precisely the sense that we are discussing for neurons.

Zipf's law is a power law, $p_r \propto r^{-\gamma}$, in this case with $\gamma = 1$, although this is quite different from the usual power law scaling relations among thermodynamic variables near an equilibrium critical point (7–9). The ubiquity of Zipf's law has led many people to wonder if there is some universal underlying mechanism. In several ways, this discussion parallels the discussion of $1/f$ noise: In many systems, fluctuations over time have a spectrum without any obvious scale, and when we plot the spectrum vs. frequency, especially on logarithmic axes, the behavior approximates a power law with exponent close to one.

In the discussion of $1/f$ noise, it was realized, early on, that a system might appear to be scale-invariant if it has a discrete set of scales spread over a sufficiently broad range. Thus, if we look at fluctuations over time, and what we see is the sum or superposition of many processes with correlation times $\tau_1, \tau_2, \tau_3, \dots$, then, if these correlation times come from a broad distribution, the net spectrum will be nearly featureless; even a handful of correlation times, with the right spread, can give a good approximation to $1/f$ noise. It seems that this is the correct description of $1/f$ noise in metals (22). Importantly, if the apparent $1/f$ noise really is a superposition of many noise sources with a range of correlation times, then, if we can perturb these time scales, we should see measurable departures from $1/f$ behavior, and this was the experimental strategy used in sorting out the behavior of current noise in metals. What this means, of course, is that this is an example of almost $1/f$ noise, and that the small deviations from truly scale-invariant behavior are crucial.

An extreme version of the mixture model for scale-invariant behavior is discussed by van Opheusden (23), who considered populations of neurons firing independently but with a distribution of mean spike probabilities. With a proper choice of this distribution, completely independent neurons can generate a good approximation to Zipf's law at fixed N . As noted in the main text, however, the actual distribution of spike probabilities that we see in the data does not have this special property. Further, unless the distribution of spike probabilities is singular, the variance of log probability across all of the states of the network will be exactly proportional to the number of neurons that we consider, and hence such models cannot explain the supralinear growth of the heat capacity in Fig. 4, which is one of the key signatures of criticality.

Aitchison et al. (24) have suggested that the original example of Zipf's law for words in English should be explained by adapting the multiple time scale idea in $1/f$ noise. The distribution of words can be thought of as a sum over contributions from several parts of speech (nouns, verbs, adjectives, etc.), and, for each part of speech, we do not see Zipf's law but rather a distribution that has a characteristic scale; the scales for different parts of speech are different, and, when we sum over all parts of speech, we see the emergence of Zipf's law. If this is correct, then, as in the case of $1/f$ noise in metals, we must conclude that Zipf's law is not exact. Further, it should be possible to modulate the characteristic scales, or the weights given to each component of the distribution, and thereby make the deviations from Zipf's law more apparent. In metals, one can do this simply by modulating the temperature (22). In language, the scales and weights for different parts of speech vary across languages, topics, and authors, so one might expect the equivalent of the temperature modulation experiment has been done, implicitly, many times, although this is not discussed in ref. 24. With modern corpora, searching more carefully for departures from Zipf's law should be straightforward. At best, however, explaining Zipf's law as a superposition over multiple parts of speech would be a demonstration that deviations from Zipf's law are important.

In connecting Zipf's law to criticality, one must keep in mind that critical phenomena exist only in the thermodynamic limit. As we have defined it, the entropy vs. energy $S(E)$ is not a smooth function in a system of finite size, because there are discrete states with particular probabilities and hence particular energies. A differentiable function $s(\epsilon)$ emerges, as in Eq. S29, only in the limit that we consider a system with many degrees of freedom. In the example of language, to make a connection to criticality thus requires more than counting words. Instead, we should imagine text segments with a length of N letters or words, and ask how the Zipf plot evolves as a function of N . The emergence of a function $s(\epsilon)$ would correspond to the plot of $(\ln p_r)/N$ vs. $(\ln r)/N$ converging to a limit as N grows, and evidence for criticality depends on the properties of this limiting function. Thus, criticality is much more than Zipf's law at fixed N .

More General Hidden Variable Models. The idea that correlations among neurons might be inherited from the visual stimulus is one possibility among many. More generally, we might ask if the pattern of correlations could be understood as the independent response of neurons to some signal that is effectively external to the network, or at least hidden from an observer who sees only the patterns of spikes and silence. To assess this possibility, it is useful to step back and think about simpler models in statistical mechanics. Almost everything that we will say in this section is well-known in the physics literature, but it seems useful to be explicit.

Consider the mean field Ising ferromagnet, in which spins $\sigma_i = \pm 1$ experience an effective magnetic field that is proportional to the average over all of the other spins in the system, so that

$$E(\{\sigma_i\}) = -\frac{J}{2N} \sum_{i \neq j} \sigma_i \sigma_j. \quad [\text{S37}]$$

Note that the sum is over all pairs, and the factor of N ensures that the energy of the system is proportional to N . The sum over all distinct pairs is missing the term $i=j$, but, because $\sigma_i^2 = 1$, we have

$$E(\{\sigma_i\}) = -\frac{J}{2N} \sum_{i,j} \sigma_i \sigma_j + \frac{J}{2} \quad [\text{S38}]$$

$$= -\frac{J}{2N} \left(\sum_i \sigma_i \right)^2 + \frac{J}{2}. \quad [\text{S39}]$$

The probability of finding the system in any particular state $\{\sigma_i\}$ is given by (choosing units where $k_B T = 1$)

$$P(\{\sigma_i\}) \equiv \frac{1}{Z} e^{-E(\{\sigma_i\})} \quad [\text{S40}]$$

$$= \frac{1}{Z} \exp \left[\frac{J}{2N} \left(\sum_i \sigma_i \right)^2 - \frac{J}{2} \right]. \quad [\text{S41}]$$

However, we can always write

$$\exp \left[\frac{A}{2} x^2 \right] = \int \frac{dh}{\sqrt{2\pi A}} \exp \left[-\frac{1}{2A} h^2 + hx \right]. \quad [\text{S42}]$$

Applying this identity to Eq. S41, we have

$$P(\{\sigma_i\}) = e^{-J/2} \frac{1}{Z} \int \frac{dh}{\sqrt{2\pi N/J}} \sum_{\{\sigma_i\}} \exp \left[-\frac{N}{2J} h^2 + h \sum_i \sigma_i \right]. \quad [\text{S43}]$$

We can think of this, more suggestively, as

$$P(\{\sigma_i\}) \propto \int dh P(h) \prod_{i=1}^N P(\sigma_i|h), \quad [\text{S44}]$$

where $P(\sigma_i|h)$ describes the response of a single spin to an external field,

$$P(\sigma_i|h) = \frac{1}{2 \cosh(h)} e^{h\sigma_i}, \quad [\text{S45}]$$

and $P(h)$ is a distribution of fields,

$$P(h) = \frac{1}{Z_h} \exp \left[-\frac{N}{2J} h^2 + N \ln \cosh(h) \right]. \quad [\text{S46}]$$

Thus, a model in which all spins interact with one another, equally, is mathematically identical to a model in which each spin responds independently to a magnetic field chosen at random.

Once we transform from interacting spins to a distribution of fields, all of the thermodynamic behavior of the system is determined by $P(h)$ (Eq. S46). We notice that if J is small (equivalently, if T is large), the distribution of fields is close to being Gaussian with standard deviation $\delta h = \sqrt{J/N}$. If J is large, the distribution $P(h)$ becomes bimodal, with peaks at $\pm h_0(J)$ whose locations do not depend on N ; this corresponds to the spontaneous magnetization of the system. Finally, at the critical value of $J = 1$, the distribution of fields is unimodal, centered at $h = 0$, but broad,

$$P_{\text{crit}}(h) \approx \exp \left[-\frac{N}{12} h^4 + \dots \right], \quad [\text{S47}]$$

so that typical fields are $\delta h \approx 1/N^{1/4}$, much larger than $\delta h \approx 1/N^{1/2}$ in the high-temperature phase. In this sense, criticality is the statement that the equivalent fields have anomalously large fluctuations (25).

We can find essentially the same equivalence in a much broader class of models. Consider a collection of spins that interact through some matrix J_{ij} , so that the energy

$$E(\{\sigma_i\}) = -\frac{1}{2} \sum_{i,j} J_{ij} \sigma_i \sigma_j. \quad [\text{S48}]$$

The Hopfield model corresponds to the choice

$$J_{ij} = \frac{J}{N} \sum_{\mu=1}^K \xi_i^\mu \xi_j^\mu, \quad [\text{S49}]$$

where there are K stored memories,

$$\xi^\mu \equiv \{\xi_1^\mu, \xi_2^\mu, \dots, \xi_N^\mu\}. \quad [\text{S50}]$$

In this case, the same arguments that lead to Eqs. S44 and S46 now give

$$P(\{\sigma_i\}) \propto \int d^k \phi P(\phi) \prod_{i=1}^N P(\sigma_i|h_i), \quad [\text{S51}]$$

where the local fields

$$h_i = \sum_{\mu=1}^K \xi_i^\mu \phi_\mu, \quad [\text{S52}]$$

and

$$P(\phi) = \frac{1}{Z_\phi} \exp \left[-\frac{N}{2J} \sum_{\mu=1}^K \phi_\mu^2 + \sum_{i=1}^N \ln \cosh \left(\sum_{\mu=1}^K \xi_i^\mu \phi_\mu \right) \right]. \quad [\text{S53}]$$

As in the mean field model, if J is small, then the effective fields have a standard deviation $\delta h \approx 1/\sqrt{N}$, and, as the system approaches the critical point, this scale becomes larger by a (fractional) power of N . This shouldn't be surprising because, with K fixed as N becomes large, the Hopfield model is a mean field model (26, 27).

If we can have K independent fields, could we have as many as there are neurons in the network? This is more subtle. If we imagine that the field acting on each neuron, as defined in Eq. S52, is built out of N independent components, so that

$$h_i = \sum_{\mu=1}^N \xi_i^\mu \phi_\mu, \quad [\text{S54}]$$

then we have to be careful to be sure that the typical field is bounded. Specifically, if all of the ϕ_μ have the same variance, $\langle \phi^2 \rangle$, then the variance of the field is

$$\langle (\delta h_i)^2 \rangle = \langle \phi^2 \rangle \sum_{\mu=1}^N (\xi_i^\mu)^2. \quad [\text{S55}]$$

Clearly, we need to have $\xi_\mu \approx 1/\sqrt{N}$ to be sure that the variance of the fields is not proportional to the size of the system. Therefore, we should write $\xi_\mu = \alpha_i^\mu / \sqrt{N}$, where α_i^μ is a number of order 1. Then the correlation between the fields acting on different neurons becomes

$$\langle \delta h_i \delta h_j \rangle = \frac{1}{N} \langle \phi^2 \rangle \sum_{\mu=1}^N \alpha_i^\mu \alpha_j^\mu. \quad [\text{S56}]$$

Now, if the influences of the different field components on the different neurons (the coefficients α_i^μ) are essentially random—e.g., some neurons are “off cells” with respect to one field and “on cells” with respect to another, with no pattern in this assignment—then the sum in Eq. S56 is of N random numbers with zero mean, and hence the typical scale for the sum is $\langle \delta h_i \delta h_j \rangle \approx 1/\sqrt{N}$. This is not at all what one expects in a critical system. The only way to escape from this conclusion is for the terms α_i^μ to have some structure, which is equivalent to fixing some correlations among the fields acting on different spins. Put another way, if the system we are studying is equivalent to one in which N spins (or neurons) are reacting independently to N distinct fields, then criticality requires some form of correlation among these fields.

The role of correlations in critical behavior is even clearer in the general case where we have an arbitrary matrix of interactions J_{ij} . Then we can write

$$P(\{\sigma_i\}) \equiv \frac{1}{Z} \exp \left[\frac{1}{2} \sum_{i,j} J_{ij} \sigma_i \sigma_j \right] \quad [\text{S57}]$$

$$= \int d^N h P(\{h_i\}) \prod_{i=1}^N P(\sigma_i | h_i), \quad [\text{S58}]$$

where the distribution of fields is given by

$$P(\{h_i\}) = \frac{1}{Z'} \exp \left[-\frac{1}{2} \sum_{i,j} K_{ij} h_i h_j + \sum_i \ln \cosh h_i \right], \quad [\text{S59}]$$

where

$$Z' = Z \sqrt{(2\pi)^N \det J}, \quad [\text{S60}]$$

and K_{ij} is the matrix inverse of J_{ij} , $K_{ij} = (J^{-1})_{ij}$. This implies that the partition function can be written as an integral over the fluctuating fields,

$$Z \propto \int d^N h \exp \left[-\frac{1}{2} \sum_{i,j} K_{ij} h_i h_j + \sum_i \ln \cosh h_i \right]. \quad [\text{S61}]$$

If we think about a family of models in which the interactions J_{ij} are scaled up and down in strength, e.g., with an effective temperature $J_{ij} \rightarrow J_{ij}/T$, then there is (often) a critical point at some value of the temperature T . What happens to the probability distribution of the equivalent fields, $P(\{h_i\})$, at this critical point? It is hard to answer this question in general, but, in the well-studied examples from statistical mechanics—where the elements of the network live on a regular lattice, and the matrix K_{ij} has a structure that depends only the distance between lattice points i and j , decaying rapidly so that the dominant terms connect near neighbors—the structure of $P(\{h_i\})$ near criticality approaches the structure of ϕ^4 field theory (9). Crucially, the variance of the field at a single point, $\langle (\delta h_i)^2 \rangle$, does not acquire any anomalous N dependence at the critical point. Instead, criticality is marked by the appearance of long-range correlations among the fields at different points, so that the sum of the fields over the entire sample (the $\mathbf{k} = 0$ Fourier component) does have a diverging variance.

To summarize, almost any model of interacting spins (or neurons) can be rewritten as a model of spins that respond independently to external signals; the thermodynamic behavior is then controlled by the distribution of these signals. If the number of signals is small compared with the number of elements in the network, which corresponds to a mean field model, then, away from criticality, the typical scale of these signals is small (e.g., $\sim 1/\sqrt{N}$), and the approach to the critical point involves this scale becoming anomalously large. In the more general case where the number of signals is comparable to the number of neurons, criticality is associated not with an anomalous scale for the fluctuations of any single signal but rather with large-scale correlations among these signals.

The correlations among neurons are described by a matrix $\chi_{ij} = \langle \sigma_i \sigma_j \rangle - \langle \sigma_i \rangle \langle \sigma_j \rangle$, and it is useful to think about the eigenvalues of this matrix. In a mean field model at criticality, or in the scenario described in ref. 25, there will be one eigenvalue separated from all of the others, which carries most of the variance of the entire system. In a system with homogeneous local interactions, the eigenmodes of χ_{ij} are Fourier modes, and, at criticality, the spectrum $\tilde{\chi}(\mathbf{k})$ diverges as $\mathbf{k} \rightarrow 0$, but continuously, so that no single mode separates cleanly from all of the others. Similarly, a mean field model is equivalent to an interacting model in which the matrix J_{ij} is of low rank (in the simplest case, rank one). Analyzing the raw data from our population of $N = 160$ neurons, we find that the largest eigenvalue of χ_{ij} captures less than 10% of the total variance, and is separated from the second largest eigenvalue by a factor of less than 2. Analyzing the models we have constructed, the spectrum of J_{ij} is nearly continuous, with no sign of a single dominant mode. These observations indicate that the network we are studying is not in the mean field regime and, more generally, that its collective behavior cannot be captured by linear dimensionality reduction strategies.

The idea that we can explain what we observe in the population of RGCs as being the result of neurons responding to other signals

evidently has clear limits. Except in special cases, assigning apparent critical behavior to such a model effectively transfers the problem to explaining the strong correlations among the signals that are driving the neurons. In this regard, it is interesting that, although details vary, we see signs of near-critical behavior in response to naturalistic movies, random checkerboards, and full-field flicker (Fig. S4). Across these different stimulus ensembles, the correlation structure of signals at the input to the retina is changing dramatically, and so, even if we think that the behavior of the ganglion cell population should be ascribed to the statistics of input signals, one has to explain how the correlations needed to mimic criticality are maintained by the retinal circuitry.

Latent Variables Redux. Aitchison et al. claim that critical phenomenology is a generic consequence of large fluctuations in latent variables (24), arguing that the behavior of mean field systems discussed by Schwab et al. (25) is typical of what we should see in complex, biological contexts. They also propose explicit candidates for the latent variables in the case of RGCs. We have argued in the preceding paragraphs that the mean field case is not the typical one in statistical physics, and is unlikely to describe the data we are discussing here. Nonetheless, one might worry that the “large variance” scenario does succeed in producing something like the critical behavior we have identified, without any of the fine tuning that one might have expected by analogy with known equilibrium critical phenomena. Here we test the suggestions of ref. 24 in more detail.

Concretely, Aitchison et al. (24) propose that all of the phenomenology of criticality should be understood in terms of models where different neurons spike or remain silent independently given the value of a latent variable that is broadcast to the entire network. Their first suggestion is that this variable is the visual stimulus itself, parameterized by time during the stimulus movie (figure 3 B and C of ref. 24). If this model is correct, then, in experiments with repeated presentations of the same stimulus, we should find zero correlations among neurons when we average over repetitions at the same moment in time. This test is not as simple as it sounds, however, because natural movies have many epochs in which the probability of one cell spiking is essentially zero. In our data, we record from $N = 160$ neurons and we have $T/\Delta\tau = 953$ distinct time bins in the natural movie, so there are $\sim 10^5$ entries in the matrix of spike probability vs. time; of these, a fraction 0.68 are consistent with zero, in that we see no spikes across 297 repetitions of the movie. Evidently, in such silent bins, one cannot estimate correlations. Conversely, a significant component of the overall correlations between two neurons may be contributed by the temporal coincidence of these silent epochs.

As an aside about silent epochs, we note that the maximum entropy model (Eqs. 3 and 4) predicts that individual neurons should have near-zero probability of spiking when the effective field contributed by the other neurons in the network is sufficiently negative. This prediction is quantitatively correct, down to probabilities of ~ 0.001 , as shown in figure 9 of ref. 11. By tracking the effective field vs. time during the stimulus movie, we can correctly predict continuous epochs of silence, characteristic of the neural response to natural stimuli, as shown in figure 15A of ref. 11.

To test the hypothesis of independence given the latent time variable, we compute the correlation coefficient between the binary variables σ_i and σ_j for every pair of cells (i, j) at a fixed time t , but only at times in the movie where each cell in the pair generates at least five spikes across the 297 repetitions of the stimulus; results are shown in Fig. S6. Although we can find moments in time where neurons that are strongly correlated across the whole experiment have near-zero correlation, we can also find the opposite. In fact, the range of correlation coefficients that we observe while conditioning on a particular time in the stimulus movie is broader than the distribution that we see in the overall correlations. It is perhaps most striking that neurons with near-zero pairwise correlation across the whole experiment can have large

positive or negative correlations when conditioned on the stimulus movie, exactly the opposite of what Aitchison et al. (24) predict.

A further difficulty in testing the hypothesis of conditional independence arises from the limited size of our data set, or any reasonable data set. As we have discussed in ref. 11, the long duration of the experiments we are analyzing means that overall correlations can be estimated with high precision, and the threshold for reliable detection of a correlation is correspondingly small. However, if we are trying to estimate the correlations at a single moment in time, even uncorrelated neurons will exhibit spurious correlations with typical scale $1/\sqrt{N_{\text{reps}}}$, where N_{reps} is the number of repetitions of the stimulus movie; even with the relatively large $N_{\text{reps}} = 297$ in this experiment, we expect spurious correlations of ~ 0.05 , as indicated in Fig. S6. The fact that many of the correlations we observe are smaller than this, of course, does not mean that neurons are conditionally independent but rather that we can't tell. This is a serious problem in drawing conclusions about the collective behavior of the network, because we know that widespread correlations on the order of $1/N$ can be signatures of nontrivial collective behavior (13, 28). To reliably exclude correlations on this scale at one moment in time, with no further assumptions, would require $N_{\text{reps}} \approx N^2$, which rapidly becomes impossible in larger networks; even if we are more optimistic and assume that the relevant scale of correlations is $\sim 1/\sqrt{N}$, we still need $N_{\text{reps}} \approx N$. This means that, with reasonable data sets on large networks, one could easily conclude that the data are statistically consistent with the hypothesis of conditional independence when, in fact, the correlations are sufficiently strong to provide the signature of dramatic collective behavior. For a different approach to this problem, which reaches similar conclusions to our Fig. S6, see ref. 29.

The second suggestion of Aitchison et al. (24) is that the relevant latent variable is the total number of spikes generated by the network, $K = \sum_{i=1}^N \sigma_i$ (figure 3 D and E of ref. 24). This is difficult to understand because K is a collective variable, not a latent variable. For a network with a finite number of neurons, for example, it is not possible for the activity of each cell to be independent given K ; at a minimum, there must be anticorrelations that hold the number of spikes fixed. In trying to make sense out of these ideas, we have examined the correlations between pairs of cells at fixed K , and find almost all possible behaviors, including strong positive correlations (the opposite of what is required to hold K fixed) with K -dependent strengths.

Our earlier work emphasized that the distribution of K itself is anomalous, and that maximum entropy models that capture this distribution already exhibit signatures of criticality (30). Focusing on the summed activity of a network of neurons is analogous to focusing on the total magnetization of a magnet. Indeed, criticality in ferromagnets is associated with an anomalously broad distribution of magnetizations, just as the signatures that we see of critical behavior in a neural network are associated with an anomalously broad distribution of K . However, in no sense does this explain the critical phenomena. In particular, the qualitative observation of large fluctuations in magnetization is consistent with many different quantitative critical behaviors, including the mean field case where there is no divergence of the specific heat.

To summarize, the suggestion by Aitchison et al. (24) that time in the stimulus movie provides a latent variable whose variation explains the behavior that we see fails because the correlations conditioned on this latent variable are as large and structured as observed without conditioning. Their suggestion that the total number of spikes is the relevant variable confuses latent with collective variables, and we find that conditioning on this collective variable also does not simplify the correlation structure of the network. We also have examples in equilibrium statistical mechanics where (qualitatively) large fluctuations in a collective variable are associated with critical phenomena of different universality classes, so that such fluctuations alone cannot single out behaviors such as those we observe in Figs. 3 and 4.

We can also assess the claim that large fluctuations in a latent variable lead generically to critical behavior by exploring a biologically plausible model. Imagine that the sensory stimulus can be parameterized by a variable \vec{x} . This could represent, in the retina, the position of a single object. More abstractly, we can think about the parameters in a space of possible stimuli, so that \vec{x} represents position in a “feature space.” We could also imagine that we are recording from a part of the brain that represents the organism’s own position in space, as with place cells in the hippocampus, in which case \vec{x} is again a literal position variable. As a model, we will consider neurons such that each cell n generates spikes when \vec{x} is in the neighborhood of that cell’s preferred stimulus \vec{x}_n . More quantitatively, we give each cell a receptive field such that the probability of spiking in a small window of time is

$$p_n(\vec{x}) = P_0 \exp\left[-\frac{|\vec{x} - \vec{x}_n|^2}{2\sigma^2}\right], \quad [\text{S62}]$$

and each cell is independent of the rest given the value of the stimulus \vec{x} . Notice that because we will be analyzing only the

distribution of responses in a single small window of time Δt , as with our analysis of the real data, we don’t need to make any assumptions about the temporal statistics of the spikes.

We focus on the simplest case, where \vec{x} is one-dimensional, and take the distribution of this variable to be uniform across some interval; without loss of generality, we can take $0 < x < 1$. We assume that the N neurons have preferred stimuli x_n that are random but uniformly distributed throughout this interval. Then the only parameters to be adjusted are the width σ of the receptive fields and the peak spike probability P_0 . Fig. S7 shows an example with $N = 100$ neurons, $P_0 = 0.3$, and $\sigma = 0.1$; reasonable variations in these parameters do not change the qualitative picture. We can generate long samples of data from this model, and then perform exactly the same analysis that we have done for the real neurons. We see that, although spike probabilities are being modulated in a correlated fashion across the entire population, there is no hint of Zipf’s law (Fig. S7B), and the plot of entropy vs. energy is far from linear (Fig. S7D). This thermodynamic signature of criticality thus is not a generic consequence of strong driving by some latent variable.

- Ruelle D (1978) *Thermodynamic Formalism: The Mathematical Structures of Classical Equilibrium Statistical Mechanics* (Addison-Wesley, Reading, MA).
- Halsey TC, Jensen MH, Kadanoff LP, Procaccia I, Shraiman BI (1986) Fractal measures and their singularities: The characterization of strange sets. *Phys Rev A* 33(2):1141–1151.
- Feigenbaum MJ, Jensen MH, Procaccia I (1986) Time ordering and the thermodynamics of strange sets: Theory and experimental tests. *Phys Rev Lett* 57(13):1503–1506.
- Stephens GJ, Mora T, Tkačik G, Bialek W (2013) Statistical thermodynamics of natural images. *Phys Rev Lett* 110(1):018701.
- Mora T, Bialek W (2011) Are biological systems poised at criticality? *J Stat Phys* 144(2): 268–302.
- TM Cover TM, Thomas JA (1991) *Elements of Information Theory* (Wiley, New York).
- Wilson KG (1979) Problems in physics with many scales of length. *Sci Am* 241(2): 158–179.
- Sethna JP (2006) *Statistical Mechanics: Entropy, Order Parameters, and Complexity* (Oxford Univ Press, Oxford, UK).
- Parisi G (1988) *Statistical Field Theory* (Addison-Wesley, Redwood City, CA).
- Guckenheimer J, Holmes P (1983) *Nonlinear Oscillations, Dynamical Systems, and Bifurcations of Vector Fields* (Springer-Verlag, New York).
- Tkačik G, et al. (2014) Searching for collective behavior in a large network of sensory neurons. *PLoS Comput Biol* 10(1):e1003408.
- Marre O, et al. (2012) Mapping a complete neural population in the retina. *J Neurosci* 32(43):14859–14873.
- Schneidman E, Berry MJ, 2nd, Segev R, Bialek W (2006) Weak pairwise correlations imply strongly correlated network states in a neural population. *Nature* 440(7087): 1007–1012.
- Mora T, Deny S, Marre O (2014) Dynamical criticality in the collective activity of a population of retinal neurons. arXiv:1410.6769 [q-bio.NC].
- Wang F, Landau DP (2001) Efficient, multiple-range random walk algorithm to calculate the density of states. *Phys Rev Lett* 86(10):2050–2053.
- Perkel DH, Bullock TH (1968) Neural coding. *Neurosci Res Prog Sum* 3:221–348.
- Mastrorade DN (1983) Interactions between ganglion cells in cat retina. *J Neurophysiol* 49(2):350–365.
- Brivanlou IH, Warland DK, Meister M (1998) Mechanisms of concerted firing among retinal ganglion cells. *Neuron* 20(3):527–539.
- Palmer SE, Marre O, Berry MJ, 2nd, Bialek W (2015) Predictive information in a sensory population. *Proc Natl Acad Sci USA* 112(22):6908–6913.
- Zipf GK (1932) *Selected Studies of the Principles of Relative Frequency in Language* (Harvard Univ Press, Cambridge, MA).
- Mora T, Walczak AM, Bialek W, Callan CG, Jr (2010) Maximum entropy models for antibody diversity. *Proc Natl Acad Sci USA* 107(12):5405–5410.
- Dutta P, Horn PM (1981) Low-frequency fluctuations in solids: $1/f$ noise. *Rev Mod Phys* 53:497–516.
- van Opheusden SCF (2013) Critical states in retinal population codes, Masters thesis (Universiteit Leiden, Leiden, The Netherlands).
- Aitchison L, Corradi N, Latham PE (2014) Zipf’s law arises naturally in structured, high-dimensional data. arXiv:1407.7135 [q-bio.NC].
- Schwab DJ, Nemenman I, Mehta P (2014) Zipf’s law and criticality in multivariate data without fine-tuning. *Phys Rev Lett* 113(6):068102.
- Amit DJ, Gutfreund H, Sompolinsky H (1987) Statistical mechanics of neural networks near saturation. *Ann Phys* 173:30–67.
- Amit DJ (1989) *Modeling Brain Function: The World of Attractor Neural Networks* (Cambridge Univ Press, Cambridge, UK).
- Castellana M, Bialek W (2014) Inverse spin glass and related maximum entropy problems. *Phys Rev Lett* 113(11):117204.
- Granot-Atedgi E, Tkačik G, Segev R, Schneidman E (2013) Stimulus-dependent maximum entropy models of neural population codes. *PLoS Comput Biol* 9(3):e1002922.
- Tkačik G, et al. (2013) The simplest maximum entropy model for collective behavior in a neural network. *J Stat Mech* 2013:P03011.

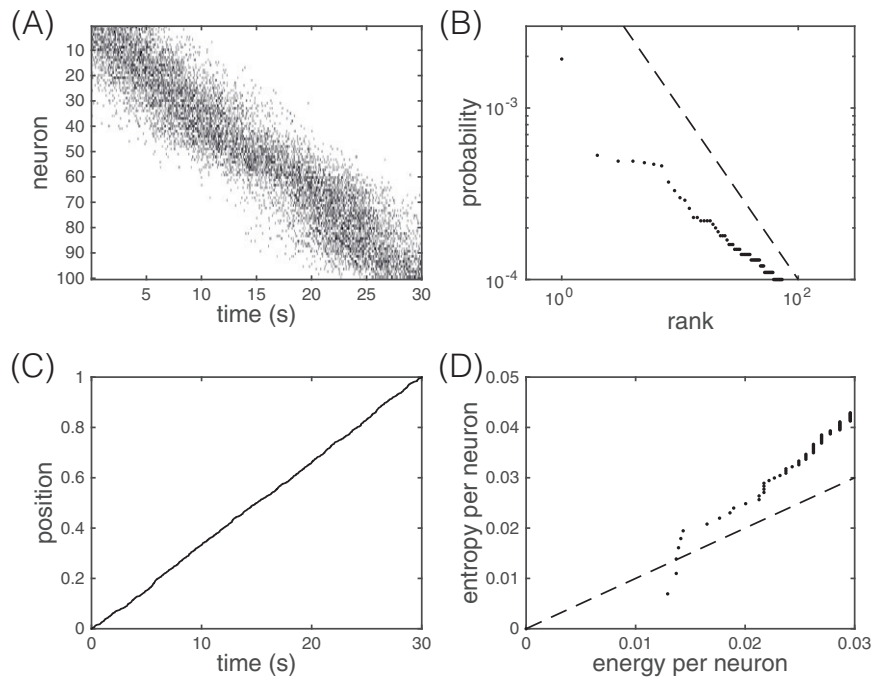


Fig. S7. Responses and thermodynamics for a population of model neurons. (A) Spike raster from a population of neurons with responses determined by Eq. S62, as the stimulus variable x moves along the trajectory shown in C. (B) Zipf plot— $\log(\text{probability})$ vs. $\log(\text{rank})$ —for the “words” describing the patterns of response in the model population of 100 cells; dashed line is Zipf’s law, for comparison. (D) Entropy vs. energy per neuron in the model population of 100 cells, computed as for the real data in Fig. 3A; dashed line is of unit slope, for comparison.