



**HAL**  
open science

# Large-Deviation Approach to Random Recurrent Neuronal Networks: Parameter Inference and Fluctuation-Induced Transitions

Alexander van Meegen, Tobias Kühn, Moritz Helias

► **To cite this version:**

Alexander van Meegen, Tobias Kühn, Moritz Helias. Large-Deviation Approach to Random Recurrent Neuronal Networks: Parameter Inference and Fluctuation-Induced Transitions. *Physical Review Letters*, 2021, 127 (15), 10.1103/physrevlett.127.158302 . hal-03407444

**HAL Id: hal-03407444**

**<https://hal.sorbonne-universite.fr/hal-03407444>**

Submitted on 28 Oct 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Large-Deviation Approach to Random Recurrent Neuronal Networks: Parameter Inference and Fluctuation-Induced Transitions


Alexander van Meegen<sup>1,2,\*</sup>, Tobias Kühn<sup>1,3,4</sup> and Moritz Helias<sup>1,3</sup>

<sup>1</sup>*Institute of Neuroscience and Medicine (INM-6) and Institute for Advanced Simulation (IAS-6) and JARA-Institute Brain Structure-Function Relationships (INM-10), Jülich Research Centre, 52428 Jülich, Germany*

<sup>2</sup>*Institute of Zoology, University of Cologne, 50674 Cologne, Germany*

<sup>3</sup>*Department of Physics, Faculty 1, RWTH Aachen University, 52074 Aachen, Germany*

<sup>4</sup>*Laboratoire de Physique de l'École Normale Supérieure, ENS, Université PSL, CNRS, Sorbonne Université, Université de Paris, F-75005 Paris, France*

 (Received 21 September 2020; revised 5 July 2021; accepted 19 August 2021; published 7 October 2021)

We here unify the field-theoretical approach to neuronal networks with large deviations theory. For a prototypical random recurrent network model with continuous-valued units, we show that the effective action is identical to the rate function and derive the latter using field theory. This rate function takes the form of a Kullback-Leibler divergence which enables data-driven inference of model parameters and calculation of fluctuations beyond mean-field theory. Lastly, we expose a regime with fluctuation-induced transitions between mean-field solutions.

DOI: [10.1103/PhysRevLett.127.158302](https://doi.org/10.1103/PhysRevLett.127.158302)

*Introduction.*—Biological neuronal networks are systems with many degrees of freedom and intriguing properties: their units are coupled in a directed, nonsymmetric manner, so that they typically operate outside thermodynamic equilibrium [1,2]. The primary analytical method to study neuronal networks has been mean-field theory [3–8]. Its field-theoretical basis has been exposed only recently [9,10]. However, to understand the parallel and distributed information processing performed by neuronal networks, the study of the forward problem—from the microscopic parameters of the model to its dynamics—is not sufficient. One additionally faces the inverse problem of determining the parameters of the model given a desired dynamics and thus function. Formally, one needs to link statistical physics with concepts from information theory and statistical inference.

We here expose a tight relation between statistical field theory of neuronal networks, large deviations theory, information theory, and inference. To this end, we generalize the probabilistic view of large deviations theory, which yields rigorous results for the leading-order behavior in the network size  $N$  [11,12], to arbitrary single unit dynamics, transfer functions, and multiple populations. We furthermore show that the central quantity of large deviations theory, the rate function, is identical to the effective action in statistical field theory. This link exposes a second

relation: Bayesian inference and prediction are naturally formulated within this framework, spanning the arc to information processing. Concretely, we develop a method for parameter inference from transient data for single- and multi-population networks. Lastly, we overcome the inherent limit of mean-field theory—its neglect of fluctuations. We develop a theory for fluctuations of the order parameter when the intrinsic timescale is large and discover a regime with fluctuation-induced transitions between two coexisting mean-field solutions.

First, we introduce the model in its most general form. Then, we develop the theory for a single population. Last, we generalize it to multiple populations.

*Model.*—We consider block-structured random networks of  $N = \sum_{\alpha} N_{\alpha}$  nonlinearly interacting units  $x_i^{\alpha}(t)$  driven by an external input  $\xi_i^{\alpha}(t)$ . The dynamics of the  $i$ th unit in the  $\alpha$ th population is governed by the stochastic differential equation

$$\tau_{\alpha} \dot{x}_i^{\alpha}(t) = -U'_{\alpha}(x_i^{\alpha}(t)) + \sum_{\beta} \sum_{j=1}^{N_{\beta}} J_{ij}^{\alpha\beta} \phi(x_j^{\beta}(t)) + \xi_i^{\alpha}(t). \quad (1)$$

In the absence of recurrent and external inputs, the units undergo an overdamped motion with time constant  $\tau_{\alpha}$  in a potential  $U_{\alpha}(x)$ . The  $J_{ij}^{\alpha\beta}$  are independent and identically Gaussian-distributed random coupling weights with zero mean and population-specific variance  $\langle (J_{ij}^{\alpha\beta})^2 \rangle = g_{\alpha\beta}^2 / N_{\beta}$  where the coupling strength  $g_{\alpha\beta}$  controls the heterogeneity of the weights. The time-varying external inputs  $\xi_i^{\alpha}(t)$  are independent Gaussian white-noise processes with

Published by the American Physical Society under the terms of the [Creative Commons Attribution 4.0 International license](https://creativecommons.org/licenses/by/4.0/). Further distribution of this work must maintain attribution to the author(s) and the published article's title, journal citation, and DOI.

zero mean and correlation functions  $\langle \xi_i^\alpha(t_1) \xi_j^\beta(t_2) \rangle = 2D_\alpha \delta_{ij} \delta_{\alpha\beta} \delta(t_1 - t_2)$ . The single-population model corresponds to the one studied in Ref. [4] if the external input vanishes,  $D = 0$ , the potential is quadratic,  $U(x) = \frac{1}{2}x^2$ , and the transfer function is sigmoidal,  $\phi(x) = \tanh(x)$ ; for  $D = \frac{1}{2}$ ,  $U(x) = -\log(A^2 - x^2)$ , and  $\phi(x) = x$  it corresponds to the one in Ref. [11], which is inspired by the dynamical spin glass model of Ref. [13].

*Field theory.*—The field-theoretical treatment of Eq. (1) employs the Martin-Siggia-Rose-de Dominicis-Janssen path integral formalism [14–17]. We denote the expectation over paths across different realizations of the noise  $\xi$  as [[18], Section A.1]

$$\langle \cdot \rangle_{x|J} \equiv \langle \langle \cdot \rangle_{x|J, \xi} \rangle_{\xi} = \int \mathcal{D}x \int \mathcal{D}\tilde{x} \cdot e^{S_0(x, \tilde{x}) - \tilde{x}^T J \phi(x)},$$

where  $\langle \cdot \rangle_{x|J, \xi}$  integrates over the unique solution of Eq. (1) given one realization  $\xi$  of the noise. Here,  $S_0(x, \tilde{x}) = \tilde{x}^T [\dot{x} + U'(x)] + D\tilde{x}^T \ddot{x}$  is the action of the uncoupled neurons. We use the shorthand notation  $\mathbf{a}^T \mathbf{b} = \sum_{i=1}^N \int_0^T dt a_i(t) b_i(t)$ .

For large  $N$ , the system becomes self-averaging, a property known from many disordered systems with large numbers of degrees of freedom: the collective behavior is stereotypical, independent of the realization  $J_{ij}$ . A self-averaging observable has a sharply peaked distribution over realizations of  $\mathbf{J}$ —the observable always attains the same value, close to its average. This, however, only holds for observables averaged over all units, reminiscent of the central limit theorem. These are generally of the form  $\sum_{i=1}^N \ell(x_i)$ , where  $\ell$  is an arbitrary functional of a single unit's trajectory. It is therefore convenient to introduce the scaled cumulant-generating functional

$$W_N(\ell) := \frac{1}{N} \ln \langle \langle e^{\sum_{i=1}^N \ell(x_i)} \rangle_{x|J} \rangle_{\mathbf{J}}, \quad (2)$$

where the prefactor  $1/N$  makes sure that  $W_N$  is an intensive quantity, reminiscent of the bulk free energy [24]. In fact, we will show that the  $N$  dependence vanishes in the limit  $N \rightarrow \infty$  because the system decouples.

Performing the average over  $\mathbf{J}$ , i.e., evaluating  $\langle e^{-\tilde{x}^T J \phi(x)} \rangle_{\mathbf{J}}$ , and introducing the auxiliary field

$$C(t_1, t_2) := \frac{1}{N} \sum_{i=1}^N \phi(x_i(t_1)) \phi(x_i(t_2)) \quad (3)$$

as well as the conjugate field  $\tilde{C}$ , we can write  $W_N$  as [ [18], Section A.1]

$$W_N(\ell) = \frac{1}{N} \ln \int \mathcal{D}C \int \mathcal{D}\tilde{C} e^{-NC^T \tilde{C} + N\Omega_\ell(C, \tilde{C})}, \quad (4)$$

$$\Omega_\ell(C, \tilde{C}) := \ln \int \mathcal{D}x \int \mathcal{D}\tilde{x} e^{S_0(x, \tilde{x}) + \frac{1}{2}\tilde{x}^T C \tilde{x} + \tilde{x}^T \tilde{C} \phi(x)}.$$

The effective action is defined as the Legendre transform of  $W_N(\ell)$ ,

$$\Gamma_N(\mu) := \int \mathcal{D}x \mu(x) \ell_\mu(x) - W_N(\ell_\mu), \quad (5)$$

where  $\ell_\mu$  is determined implicitly by the condition  $\mu = W'_N(\ell_\mu)$  and the derivative  $W'_N(\ell)$  has to be understood as a generalized derivative, the coefficient of the linearization akin to a Fréchet derivative [25].

Note that  $W_N$  and  $\Gamma_N$  are, respectively, generalizations of a cumulant-generating functional and of the effective action [26] because both map a functional ( $\ell$  or  $\mu$ ) to the reals. For the choice  $\ell(x) = j^T x$ , where  $j(t)$  is an arbitrary function, we recover the usual cumulant-generating functional of the single unit's trajectory [ [18], Section A.4] and the corresponding effective action.

*Rate function.*—Any network-averaged observable, for which we may expect self-averaging to hold, can likewise be obtained from the empirical measure

$$\mu(y) := \frac{1}{N} \sum_{i=1}^N \delta(x_i - y), \quad (6)$$

since  $(1/N) \sum_{i=1}^N \ell(x_i) = \int \mathcal{D}y \mu(y) \ell(y)$ . Of particular interest is the leading-order exponential behavior of the distribution of empirical measures  $P(\mu) = \langle \langle P(\mu|\mathbf{x}) \rangle_{x|J} \rangle_{\mathbf{J}}$  across realizations of  $\mathbf{J}$  and  $\xi$ . This behavior in the large  $N$  limit is described by what is known as the rate function

$$H(\mu) := -\lim_{N \rightarrow \infty} \frac{1}{N} \ln P(\mu) \quad (7)$$

in large deviations theory [see, e.g., [27]];  $H(\mu)$  captures the leading exponential probability  $P(\mu) \stackrel{N \gg 1}{\simeq} \exp[-NH(\mu)]$ . For large  $N$ , the probability of an empirical measure that does not correspond to the minimum  $H'(\bar{\mu}) = 0$  is thus exponentially suppressed. Put differently, the system is self-averaging and the statistics of any network-averaged observable can be obtained using  $\bar{\mu}$ .

Similar as in field theory, it is convenient to introduce the scaled cumulant-generating functional of the empirical measure. Because  $(1/N) \sum_{i=1}^N \ell(x_i) = \int \mathcal{D}y \mu(y) \ell(y)$  holds for an arbitrary functional  $\ell(x_i)$  of the single unit's trajectory  $x_i$ , Eq. (2) has the form of the scaled cumulant-generating functional for  $\mu$  at finite  $N$ .

Using a saddle-point approximation for the integrals over  $C$  and  $\tilde{C}$  in Eq. (4) [ [18], Section A.1], we get

$$W_\infty(\ell) = -C_\ell^T \tilde{C}_\ell + \Omega_\ell(C_\ell, \tilde{C}_\ell). \quad (8)$$

Both  $C_\ell$  and  $\tilde{C}_\ell$  are determined self-consistently by the saddle-point equations  $C_\ell = \partial_{\tilde{C}} \Omega_\ell(C, \tilde{C})|_{C_\ell, \tilde{C}_\ell}$  and

$\tilde{C}_\ell = \partial_C \Omega_\ell(C, \tilde{C})|_{C_\ell, \tilde{C}_\ell}$  where  $\partial_C$  denotes a partial functional derivative.

From the scaled cumulant-generating functional, Eq. (8), we obtain the rate function via a Legendre transformation [28]:  $H(\mu) = \int \mathcal{D}x \mu(x) \ell_\mu(x) - W_\infty(\ell)$  with  $\ell_\mu$  implicitly defined by  $\mu = W'_\infty(\ell_\mu)$ . Note that  $H(\mu)$  is still convex even if  $\mu$  itself is multimodal. Comparing with Eq. (5), we observe that the rate function is equivalent to the effective action:  $H(\mu) = \lim_{N \rightarrow \infty} \Gamma_N(\mu)$ . The equation  $\mu = W'_\infty(\ell_\mu)$  can be solved for  $\ell_\mu$  to obtain a closed expression for the rate function viz. effective action [ [18], Section A.2], one main result of our work,

$$H(\mu) = \int \mathcal{D}x \mu(x) \ln \frac{\mu(x)}{\langle \delta(\dot{x} + U'(x) - \eta) \rangle_\eta}, \quad (9)$$

where  $\eta$  is a zero-mean Gaussian process with a correlation function that is determined by  $\mu(x)$ ,

$$C_\eta(t_1, t_2) = 2D\delta(t_1 - t_2) + g^2 \int \mathcal{D}x \mu(x) \phi(x(t_1)) \phi(x(t_2)). \quad (10)$$

For  $D = \frac{1}{2}$ ,  $U(x) = -\log(A^2 - x^2)$ , and  $\phi(x) = x$ , Eq. (9) can be shown to be equivalent to the mathematically rigorous result obtained in the seminal work by Ben Arous and Guionnet [ [18], Section A.3].

The rate function Eq. (9) takes the form of a Kullback-Leibler divergence. Thus, it possesses a minimum at

$$\bar{\mu}(x) = \langle \delta(\dot{x} + U'(x) - \eta) \rangle_\eta. \quad (11)$$

This most likely measure corresponds to the well-known self-consistent stochastic dynamics that is obtained in field theory [4,9,10,29]. Note that the correlation function of the effective stochastic input  $\eta$  at the minimum depends self-consistently on  $\bar{\mu}(x)$  through Eq. (10). However, the rate function  $H(\mu)$  contains more information. It quantifies the suppression of departures  $\mu - \bar{\mu}$  from the most likely measure and therefore allows the assessment of fluctuations that are beyond the scope of the classical mean-field result.

*Parameter inference.*—The rate function opens the way to address the inverse problem: given the network-averaged activity statistics, encoded in the corresponding empirical measure  $\mu$ , what are the statistics of the connectivity and the external input, i.e.,  $g$  and  $D$ ?

We determine the parameters using maximum likelihood estimation. Using Eq. (7) and Eq. (9), the likelihood of the parameters is given by

$$\ln P(\mu|g, D) \simeq -NH(\mu|g, D),$$

where  $\simeq$  denotes equality in the limit  $N \rightarrow \infty$  and we made the dependence on  $g$  and  $D$  explicit. The maximum

likelihood estimate of the parameters  $g$  and  $D$  corresponds to the minimum of the Kullback-Leibler divergence  $H$ , Eq. (9), on the right-hand side. Evaluating the derivative of  $H(\mu|g, D)$  yields [ [18], Section B.1]

$$\partial_a \ln P(\mu|g, D) \simeq -\frac{N}{2} \text{tr} \left( (C_0 - C_\eta) \frac{\partial C_\eta^{-1}}{\partial a} \right),$$

where we abbreviated  $a \in \{g, D\}$  and defined  $C_0(t_1, t_2) \equiv \int \mathcal{D}x \mu(x) (\dot{x}(t_1) + U'(x(t_1))) (\dot{x}(t_2) + U'(x(t_2)))$ . The derivative vanishes for  $C_0 = C_\eta$ . Assuming stationarity, in the Fourier domain this condition reads

$$\mathcal{S}_{\dot{x}+U'(x)}(f) = 2D + g^2 \mathcal{S}_{\phi(x)}(f), \quad (12)$$

where  $\mathcal{S}_X(f)$  denotes the network-averaged power spectrum of the observable  $X$ . Using non-negative least squares [30], Eq. (12) allows a straightforward inference of  $g$  and  $D$  (Fig. 1). To determine the transfer function  $\phi$  and the potential  $U$ , one can use model comparison techniques [ [18], Section B.2]. Using the inferred parameters, we can also predict the future activity of a unit from the knowledge of its recent past [ [18], Section B.3].

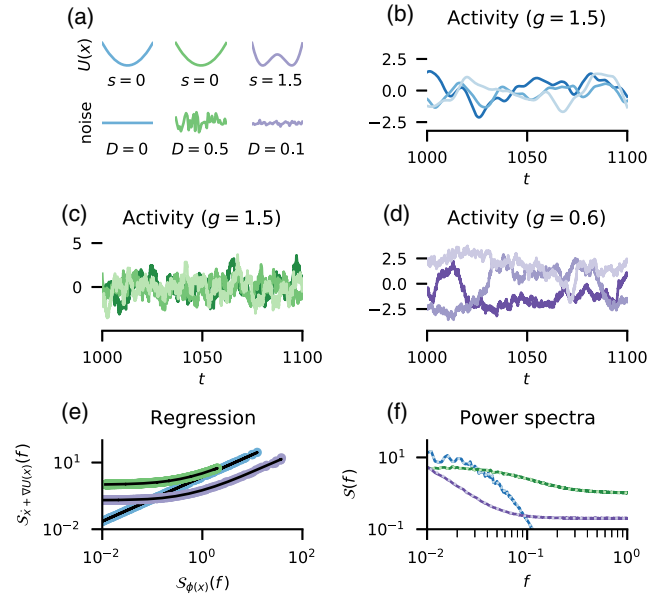


FIG. 1. Maximum likelihood parameter estimation for  $\phi(x) = \text{erf}(\sqrt{\pi}x/2)$ , potential  $U(x) = \frac{1}{2}x^2 + s \ln \cosh x$ , and external noise  $D$ . (a) Color-coded sketch of potential and noise. (b)–(d) Activity of three randomly chosen units for coupling strengths  $g$  indicated in title. (e) Parameter estimation via non-negative least squares regression (black lines) based on Eq. (12). (f) Power spectra on the left- (dark, solid curves) and right-hand sides (light, dotted curves) of Eq. (12) for the inferred parameters. Further parameters:  $\tau = 1$ ,  $N = 10000$ , temporal discretization  $dt = 10^{-2}$ , simulation time  $T = 1000$ , time span discarded to reach steady state  $T_0 = 100$ .

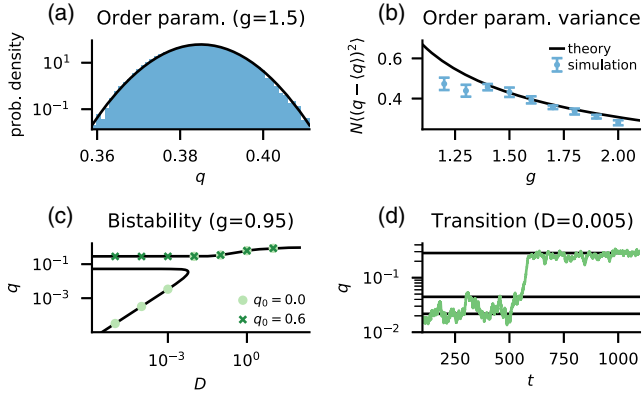


FIG. 2. Order parameter fluctuations for  $\phi(x) = \text{erf}(\sqrt{\pi}x/2)$  [(a),(b)] and metastability for  $\phi(x) = \text{clip}[\tan(x), -1, 1]$  [(c),(d)]. (a) Temporal order parameter statistics across ten simulations (bars) and theory (solid curve) from Eq. (13). (b) Order parameter variance for 10 realizations of the connectivity with standard error of the mean (symbols) and theory (solid curve) from Eq. (13). (c) Mean order parameter for different initial values  $q_0$  from simulations (symbols) and self-consistent theory (solid curves). (d) Fluctuation-induced bistability of the order parameter for  $N = 750$ ,  $g = 0.95$ . Further parameters:  $T = 5000$  in (a),(d);  $U(x) = \frac{1}{2}x^2$  in (a)–(d); other parameters as in Fig. 1.

*Fluctuations.*—The rate function allows us to go beyond mean-field theory and examine fluctuations of the order parameter. Here, we use the network-averaged variance  $q(t) = C(t, t)$  from Eq. (3) as an order parameter and restrict the discussion to the case  $U(x) = \frac{1}{2}x^2$ .

Figure 2(a) shows the distribution of  $q(t)$  across time and across realizations of the connectivity. The fluctuations across realizations of the connectivity can be computed from the curvature of the rate function  $I(C)$  that is obtained from (9) by the contraction principle [ [18], Section C.1]. In a stationary state and considering only the fluctuations across realizations of the connectivity, for slow recurrent dynamics  $\tau_c \gg 1$  we obtain the approximation for the fluctuations of  $q$

$$\langle (q - \langle q \rangle_J)^2 \rangle_J = \frac{\langle (\phi\phi - \langle \phi\phi \rangle_0)^2 \rangle_0}{N[1 - g^2(\langle \phi''\phi \rangle_0 + \langle \phi'\phi' \rangle_0)]^2}. \quad (13)$$

Here,  $\langle fg \rangle_0 \equiv \langle f(x(t))g(x(t)) \rangle_0$  denotes an expectation with respect to the self-consistent measure (11). For vanishing noise,  $D = 0$ , and  $g > 1$ , the dynamics are slow and the theory matches the empirical fluctuations very well [Figs. 2(a) and 2(b)]. Deviations in Fig. 2(b) are caused by two effects: For  $g \searrow 1$ , periodic solutions appear as a finite-size effect; for growing  $g$ , the timescale  $\tau_c$  decreases, eventually violating the assumption  $\tau_c \gg 1$  entering Eq. (13). Rate functions like  $I(C)$  in general also allow one to estimate the tail probability  $\mathbb{P}(q > \theta) \approx \exp[-NI(\theta)]$ , which here shows a quadratic decline for large departures [Fig. 2(a)].

When the denominator in Eq. (13) vanishes, fluctuations grow large, indicative of a continuous phase transition. For  $\phi'''(0) < 0$  the denominator vanishes for  $g \geq 1$  [Fig. 2(b)], in line with the established theory, the breakdown of linear stability of the fixed point  $x = 0$  [4]. For  $\phi'''(0) > 0$ , however, Eq. (13) predicts qualitatively different behavior: the denominator vanishes at  $g < 1$ , in the linearly stable regime. In fact, we find that this regime features the coexistence of two stable mean-field solutions (Fig. 2(c), [ [18], Section C.2]) and fluctuation-driven first-order transitions between them [Fig. 2(d)]. The solution with larger  $q$  corresponds to self-sustained activity; the solution with smaller  $q$  corresponds to the fixed point  $x = 0$  and is stable [ [18], Section C.2], in contrast to the case of a threshold-power-law transfer function [6].

*Multiple populations.*—For multiple populations, any population-averaged observable can be obtained from the empirical measure  $\mu^\alpha(y) = (1/N_\alpha) \sum_{i=1}^{N_\alpha} \delta(x_i^\alpha - y)$ . The joint distribution of all population-specific empirical measures  $\{\mu^\alpha\}$  is determined by the rate function [ [18], Section D]

$$H(\{\mu^\alpha\}) = \sum_\alpha \gamma_\alpha \int \mathcal{D}x \mu^\alpha(x) \ln \frac{\mu^\alpha(x)}{\langle \delta(\tau_\alpha \dot{x} + U'_\alpha(x) - \eta_\alpha) \rangle_{\eta_\alpha}}, \quad (14)$$

where  $\gamma_\alpha = N_\alpha/N$  and  $\eta_\alpha$  is a zero-mean Gaussian process with

$$C_\eta^\alpha(t_1, t_2) = 2D_\alpha \delta(t_1 - t_2) + \sum_\beta g_{\alpha\beta}^2 \int \mathcal{D}x \mu^\beta(x) \phi(x(t_1)) \phi(x(t_2)). \quad (15)$$

Again, the rate function can be interpreted as a log-likelihood; its derivative leads to [ [18], Section E.1]

$$S_{\tau_\alpha \dot{x} + U'_\alpha(x)}^\alpha(f) = 2D_\alpha + \sum_\beta g_{\alpha\beta}^2 S_{\phi(x)}^\beta(f), \quad (16)$$

which generalizes Eq. (12) to multiple populations.

Using Eq. (16), the inferred connectivity  $g_{\alpha\beta}$  matches the ground truth well; accordingly, two unconnected populations [Figs. 3(a) and 3(b)] can be clearly distinguished from a more involved network where one population ( $\alpha = 1$ ) is only active due to the recurrent input from the other population [ $\alpha = 2$ , Figs. 3(c) and 3(d)]. The method can thus distinguish intrinsically generated activity from a case where activity is driven from outside the network. However, inference of a unique set of parameters is only possible if the output spectra  $S_{\phi(x)}^\alpha(f)$  differ sufficiently across  $\alpha$ . If the output spectra match closely, Eq. (16) leads to a degenerate set of solutions that satisfy  $\sum_\beta g_{\alpha\beta}^2 = \text{const}$  and are all equally likely given the data [ [18], Section E.2].

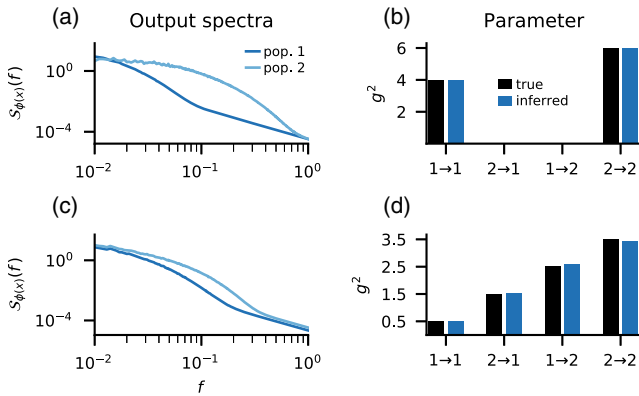


FIG. 3. Maximum likelihood parameter estimation for two populations with different time constants  $\tau_1 = 5$ ,  $\tau_2 = 1$ . (a) Output power spectra  $S_{\phi(x)}^\alpha(f)$  of two unconnected populations  $g_{21}^2 = g_{21}^2 = 0$  with  $g_{11}^2 = 4$  and  $g_{22}^2 = 6$ . (b) Estimated (blue) and true (black) parameters corresponding to (a). (c) Output power spectra of two connected populations with  $g_{11}^2 = 0.5$ ,  $g_{12}^2 = 1.5$ ,  $g_{21}^2 = 2.5$ , and  $g_{22}^2 = 3.5$ . (d) Estimated (blue) and true (black) parameters corresponding to (c). Further parameters:  $N_1 = N_2 = 5000$ ,  $\phi(x) = \text{erf}(\sqrt{\pi}x/2)$ ,  $U(x) = \frac{1}{2}x^2$ , and  $D = 0$ ; simulation parameters as in Fig. 1.

*Discussion.*—In this Letter, we found a tight link between the field-theoretical approach to neuronal networks and its counterpart based on large deviations theory. We obtained the rate function of the empirical measure for the widely used and analytically solvable model of a recurrent neuronal network [4] by field-theoretical methods. This rate function generalizes the seminal result by Ben Arous and Guionnet [11,12] to arbitrary potentials, transfer functions, and multiple populations. Intriguingly, our derivation elucidates that the rate function is identical to the effective action and takes the form of a Kullback-Leibler divergence, akin to Sanov’s theorem for sums of i.i.d. random variables [27,28]. The rate function can thus be interpreted as a distance between an empirical measure, for example given by data, and the activity statistics of the network model. This result allows us to address the inverse problem of inferring the parameters of the connectivity and external input from a set of trajectories and to determine the potential and the transfer function.

We here restricted the analysis to networks with independently drawn random weights with zero mean. Since correlated weights have a profound impact on the dynamics that can be captured using both field theory [31] and large deviations theory [32,33], it is an interesting challenge to extend the analysis in this direction. Likewise, synaptic weights with nonvanishing mean, as they appear in sparsely connected networks, present an interesting extension, because they promote fluctuation-driven states when feedback is sufficiently positive. Motifs are another important deviation from independent weights in biological neural networks are motifs [34], which pose a significant

challenge already for the field-theoretical approach [35]. Beyond the weight statistics, we assumed that the dynamics are governed by the first-order differential equation (1). Indeed, the field-theoretical approach can be generalized to a much broader class of dynamics that do not necessarily possess an action [36]; hence, it seems possible to also derive large deviations results for more general dynamics. In this regard, the extension to spiking networks is a particularly interesting but also challenging future direction. Whether the model, Eq. (1), with its current limitations—the independent weights and the first-order dynamics—allows accurate inference of network parameters from cortical recordings is an intriguing question for further research.

The unified description of random networks by statistical field theory and large deviations theory opens the door to established techniques from either domain to capture beyond mean-field behavior. Such corrections are important for small or sparse networks with nonvanishing mean connectivity, to explain correlated neuronal activity, and to study information processing in finite-size networks with realistically limited resources. We here make a first step by computing fluctuation corrections from the rate function. The quantitative theory explains near-critical fluctuations for  $g \in [1, 1 + \delta(N)]$  and we discover that expansive gain functions, as found in biology [37], lead to qualitatively different collective behavior than the well-studied contractive sigmoidal ones: The former feature metastable network states with noise-induced first order transitions between them; the latter allow for only a single solution and show second order phase transitions.

We are grateful to Olivier Faugeras and Etienne Tanré for helpful discussions on LDT of neuronal networks, to Anno Kurth for pointing us to the Fréchet derivative, and to Alexandre René, David Dahmen, Kirsten Fischer, and Christian Keup for feedback on an earlier version of the manuscript. This work was partly supported by the Helmholtz young investigator’s group VH-NG-1028, European Union Horizon 2020 Grant No. 785907 (Human Brain Project SGA2), the Human Frontier Science Program RGP0057/2016 grant, BMBF Grant “Renormalized Flows” (01IS19077A), and the Excellence Initiative of the German federal and state governments (G:(DE-82)EXS-PF-JARASDS005).

\*Corresponding author.

avm@physik.huberlin.de

- [1] M. I. Rabinovich, P. Varona, A. I. Selverston, and H. D. I. Abarbanel, *Rev. Mod. Phys.* **78**, 1213 (2006).
- [2] H. Sompolinsky, *Phys. Today* **41**, No. 12, 70 (1988).
- [3] S.-I. Amari, *IEEE Trans. SMC-2*, 643 (1972).
- [4] H. Sompolinsky, A. Crisanti, and H. J. Sommers, *Phys. Rev. Lett.* **61**, 259 (1988).
- [5] M. Stern, H. Sompolinsky, and L. F. Abbott, *Phys. Rev. E* **90**, 062710 (2014).

- [6] J. Kadmon and H. Sompolinsky, *Phys. Rev. X* **5**, 041030 (2015).
- [7] J. Aljadeff, M. Stern, and T. Sharpee, *Phys. Rev. Lett.* **114**, 088101 (2015).
- [8] A. van Meegen and B. Lindner, *Phys. Rev. Lett.* **121**, 258302 (2018).
- [9] A. Crisanti and H. Sompolinsky, *Phys. Rev. E* **98**, 062120 (2018).
- [10] J. Schuecker, S. Goedeke, and M. Helias, *Phys. Rev. X* **8**, 041029 (2018).
- [11] G. B. Arous and A. Guionnet, *Probab. Theory Relat. Fields* **102**, 455 (1995).
- [12] A. Guionnet, *Probab. Theory Relat. Fields* **109**, 183 (1997).
- [13] H. Sompolinsky and A. Zippelius, *Phys. Rev. Lett.* **47**, 359 (1981).
- [14] P. Martin, E. Siggia, and H. Rose, *Phys. Rev. A* **8**, 423 (1973).
- [15] H.-K. Janssen, *Z. Phys. B* **23**, 377 (1976).
- [16] C. Chow and M. Buice, *J. Math. Neurosci.* **5**, 8 (2015).
- [17] J. A. Hertz, Y. Roudi, and P. Sollich, *J. Phys. A* **50**, 033001 (2017).
- [18] See Supplemental Material at <http://link.aps.org/supplemental/10.1103/PhysRevLett.127.158302> for detailed derivations and further information, which includes Refs. [19–23].
- [19] J. Stapmanns, T. Kühn, D. Dahmen, T. Luu, C. Honerkamp, and M. Helias, *Phys. Rev. E* **101**, 042124 (2020).
- [20] D. J. MacKay, *Information Theory, Inference and Learning Algorithms* (Cambridge University Press, Cambridge, England, 2003).
- [21] G. Matheron, *Econ. Geol.* **58**, 1246 (1963).
- [22] C. K. Williams, *Neural Comput.* **10**, 1203 (1998).
- [23] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright *et al.*, *Nat. Methods* **17**, 261 (2020).
- [24] N. Goldenfeld, *Lectures on Phase Transitions and the Renormalization Group* (Perseus Books, Reading, Massachusetts, 1992).
- [25] M. S. Berger, *Nonlinearity and Functional Analysis*, 1st ed. (Elsevier, New York, 1977), ISBN 9780120903504.
- [26] J. Zinn-Justin, *Quantum Field Theory and Critical Phenomena* (Clarendon Press, Oxford, 1996).
- [27] M. Mezard and A. Montanari, *Information, Physics and Computation* (Oxford University Press, New York, 2009).
- [28] H. Touchette, *Phys. Rep.* **478**, 1 (2009).
- [29] M. Helias and D. Dahmen, *Statistical Field Theory for Neural Networks*, vol. 970 (Springer International Publishing, Cham, 2020).
- [30] C. L. Lawson and R. J. Hanson, *Solving Least Squares Problems* (SIAM, Philadelphia, 1995).
- [31] D. Martí, N. Brunel, and S. Ostojic, *Phys. Rev. E* **97**, 062314 (2018).
- [32] O. Faugeras and J. MacLaurin, *Entropy* **17**, 4701 (2015).
- [33] O. Faugeras, J. MacLaurin, and E. Tanré, [arXiv:1901.10248](https://arxiv.org/abs/1901.10248).
- [34] S. Song, P. Sjöström, M. Reigl, S. Nelson, and D. Chklovskii, *PLoS Biol.* **3**, e350 (2005).
- [35] D. Dahmen, S. Recanatesi, G. K. Ocker, X. Jia, M. Helias, and E. Shea-Brown, [bioRxiv](https://arxiv.org/abs/2020.021064) (2020).
- [36] C. Keup, T. Kühn, D. Dahmen, and M. Helias, *Phys. Rev. X* **11**, 021064 (2021).
- [37] A. Roxin, N. Brunel, D. Hansel, G. Mongillo, and C. van Vreeswijk, *J. Neurosci.* **31**, 16217 (2011).