

Entropy production and Kullback-Leibler divergence between stationary trajectories of discrete systems

Édgar Roldán and Juan M. R. Parrondo

Departamento de Física Atómica, Molecular y Nuclear and GISC, Universidad Complutense de Madrid, 28040 Madrid, Spain

(Received 30 January 2012; published 21 March 2012)

The irreversibility of a stationary time series can be quantified using the Kullback-Leibler divergence (KLD) between the probability of observing the series and the probability of observing the time-reversed series. Moreover, this KLD is a tool to estimate entropy production from stationary trajectories since it gives a lower bound to the entropy production of the physical process generating the series. In this paper we introduce analytical and numerical techniques to estimate the KLD between time series generated by several stochastic dynamics with a finite number of states. We examine the accuracy of our estimators for a specific example, a discrete flashing ratchet, and investigate how close the KLD is to the entropy production depending on the number of degrees of freedom of the system that are sampled in the trajectories.

DOI: [10.1103/PhysRevE.85.031129](https://doi.org/10.1103/PhysRevE.85.031129)

PACS number(s): 05.70.Ln, 05.20.-y, 05.40.-a

I. INTRODUCTION

The relationship between irreversibility and entropy production is mentioned in many undergraduate courses of thermodynamics and statistical physics. A canonical example is a glass falling to the ground and smashing into pieces. The time reverse of this process is compatible with Newton's laws, but the chances for it to occur spontaneously are incredibly small. Such a process is irreversible and the signature of this irreversibility is the production of a macroscopic amount of entropy in the universe.

The relation between irreversibility and entropy production was only a qualitative statement until the recent introduction of the Kullback-Leibler divergence (KLD) in the context of fluctuation theorems [1,2]. The time irreversibility of a process is given by the distinguishability between the process and its time reversal, which in turn can be quantified using the KLD or relative entropy, a measure of the distinguishability between two probability distributions defined in information theory [2–4]. This KLD, multiplied by the Boltzmann constant, turns out to be a lower bound to the entropy production along the process [1,2,5–15]. The bound becomes more accurate when the observables that are used to calculate the KLD contain a more complete description of the state of the system. This result has been derived in a variety of situations such as driven systems under Hamiltonian [1,2,8] and Langevin [11,14,16,17] dynamics; Markovian processes [7,10]; and also for electrical circuits [12]. Andrieux *et al.* have verified it experimentally using the data of the position of a Brownian particle in a moving optical trap [14,18], and we have shown that the bound yields a useful estimate of the entropy production in a nonequilibrium stationary state (NESS).

Imagine repeatedly sampling (or measuring) an observable of a system in a NESS. The trajectory of the outcomes is a stationary time series that can be used to estimate the KLD, by comparing the statistics of the time series with the statistics of the same series but time reversed [9]. This means that one can bound from below the entropy production in the NESS from a single time series obtained in an experiment. Such a tool is of interest in many practical situations. For instance, it allows one to discriminate between active and passive

processes in biological systems, or even to estimate or bound the amount of entropy produced, and therefore the amount of ATP consumed in a biological process. In fact, there have been previous attempts to make this distinction. Martin *et al.* have considered the violation of the fluctuation-dissipation relationship as a signature of nonequilibrium in the motion of a hair cell by using two types of measurement: the spontaneous motion of the hair bundle and the response to an external force [19]. Amman *et al.* discriminated between equilibrium and the NESS in a three-state chemical system [20]. Finally, Kennel introduced in Ref. [21] criteria based on compression algorithms to distinguish between time-symmetric and time-asymmetric chaotic series but without any connection to the physical entropy.

We are interested in estimating the KLD between the probability of observing a stationary trajectory of one or several observables of the system and the probability of observing the same trajectory but time reversed. We want to explore how this quantity bounds the entropy production of the underlying physical process [1,2,9] depending on the number of degrees of freedom of the system that are sampled in the observed stationary trajectory. Two distinct issues immediately arise: the estimation of the KLD from an empirical stationary time series and the accuracy of the bound. In this paper we address these two issues by introducing numerical and semianalytical techniques to estimate the KLD from data obtained from systems with a finite number of states.

There have been different attempts to provide accurate estimators of the KLD from a finite number of data. References [22,23] investigate how this measure can be estimated when considering empirical probability distributions of two different Markovian and higher order Markovian time series. They develop techniques based on empirical counting of finite sequences of data which are generalized to real-valued time series in Refs. [14,24,25]. A different approach is given in Ref. [26], where the KLD between two different probability distributions is estimated using compression algorithms. In this paper we refine these methods and test their performance when used to estimate the KLD from single stationary trajectories.

To explore the bound to the entropy production, we work with a discrete flashing ratchet model, where we can compare

the entropy production with the analytical value and the empirical estimations of the KLD. With this model, we can analyze how information losses affect the estimation of the KLD and the tightness of the bound for the entropy production.

The paper is organized as follows: Section II reviews the concept of the KLD and discusses its connection with entropy production. In Sec. III we present analytical and semianalytical tools to calculate the KLD between hidden Markov chains. Section IV gives a detailed description of the estimators of the KLD from empirical data, whose performance for the flashing ratchet is analyzed in Sec. V. Finally, we present our main conclusions in Sec. VI.

II. KULLBACK-LEIBLER DIVERGENCE, IRREVERSIBILITY, AND ENTROPY PRODUCTION

A. The Kullback-Leibler divergence

The Kullback-Leibler divergence, or relative entropy, measures the distinguishability of two probability distributions $p(x)$ and $q(x)$:

$$D[p(x)||q(x)] = \int dx p(x) \ln \frac{p(x)}{q(x)}. \quad (1)$$

It is always positive and vanishes if and only if $p(x) = q(x)$ for all x . Its interpretation as a measure of distinguishability is a consequence of the Chernoff-Stein lemma [3]: The probability of incorrectly guessing (via hypothesis testing) that a sequence of n data is distributed according to p when the true distribution is q is asymptotically equal to $e^{-nD[p(x)||q(x)]}$. Therefore, when p and q are similar—in the sense that they overlap significantly—the likelihood of incorrectly guessing the distribution, p or q , is large [3].

Let us recall a property of the KLD that we will use throughout the paper [3]. If we have two random variables X, Y and two joint probability distributions $p(x, y)$ and $q(x, y)$, then

$$D[p(x, y)||q(x, y)] \geq D[p(x)||q(x)]. \quad (2)$$

This means that it is harder to distinguish between p and q when we consider only the marginal distributions, $p(x)$ and $q(x)$, instead of the full joint distributions, $p(x, y)$ and $q(x, y)$. If X, Y describe the state of a system, Eq. (2) indicates that the KLD decreases when only a partial description of the system, given by the variable X , is available.

B. Irreversibility and entropy production

Consider a physical system with Hamiltonian $H(z; \lambda)$, where z denotes a point in phase space Γ and λ is a parameter of the system controlled by an external agent. The system is initially isolated in equilibrium at temperature T , and the external agent modifies λ following a protocol λ_t , with $t \in [0, \tau]$. We then let the system equilibrate by coupling it to a bath at temperature T' . The initial and final states of this process are equilibrium states for which entropy is well defined. We denote by $\rho(z, t)$ the probability density on phase space at time t , and by $\tilde{\rho}(\tilde{z}, t)$ the probability density when the system is driven by the time-reversed protocol $\tilde{\lambda}_t = \lambda_{\tau-t}$ with $t \in [0, \tau]$. Here \tilde{z} denotes the point in phase space resulting from changing the sign of all momenta in z . In Ref. [2] it is

proved that the change of the entropy ΔS in the system plus the bath, averaged over many realizations of the process, satisfies

$$\langle \Delta S \rangle = kD[\rho(z, t)||\tilde{\rho}(\tilde{z}, \tau - t)], \quad (3)$$

where k is Boltzmann's constant. Equation (3) is valid for a variety of initial equilibrium conditions [2]: canonical, multicanonical (several uncoupled systems at different temperatures), and grand-canonical distributions, as well as for different types of baths equilibrating the system at the end of the process. In particular, for canonical initial conditions in the forward and in the time-reversed processes, both at the same temperature T , Eq. (3) reads (see Ref. [2])

$$\begin{aligned} \langle \Delta S \rangle &= \langle \Delta S_{\text{system}} \rangle + \langle \Delta S_{\text{bath}} \rangle = \frac{\langle \Delta E \rangle - \Delta F}{T} + \frac{\langle Q \rangle}{T} \\ &= \frac{\langle W \rangle - \Delta F}{T}, \end{aligned} \quad (4)$$

where $\langle \Delta E \rangle$ and ΔF refer respectively to the system average energy and free energy change, Q is the heat exchanged with the thermal bath at the end of the process (realization dependent), and $W = \Delta E + Q$ is the work performed by the external agent. Therefore, in this specific case, entropy production equals the average dissipated work $\langle W_{\text{diss}} \rangle = \langle W \rangle - \Delta F$ divided by the temperature T and (3) becomes

$$\langle W_{\text{diss}} \rangle = kTD[\rho(z, t)||\tilde{\rho}(\tilde{z}, \tau - t)]. \quad (5)$$

Since the evolution is deterministic, except for the last stage where the system is connected to the bath, the point z at time t determines the whole trajectory of the system $\{z(t)\}_{t=0}^{\tau}$. Then $z(t)$ and $\{z(t)\}_{t=0}^{\tau}$ carry the same information and the KLD of their respective PDFs are equal. Equation (5) can be rewritten in terms of path probabilities \mathcal{P} [16]:

$$\langle W_{\text{diss}} \rangle = kTD[\mathcal{P}(\{z(t)\}_{t=0}^{\tau})||\tilde{\mathcal{P}}(\{\tilde{z}(\tau - t)\}_{t=0}^{\tau})]. \quad (6)$$

On the other hand, integrating Crook's relationship [27], $W - \Delta F = \ln \frac{p(W)}{\tilde{p}(-W)}$, where $p(W)$ [$\tilde{p}(-W)$] is the probability density of the work done on the system along the actual (time-reversed) process [16,27], one immediately gets

$$\langle W_{\text{diss}} \rangle = kTD[p(W)||\tilde{p}(-W)]. \quad (7)$$

Notice that the work W is a function of the trajectory $\{z(t)\}_{t=0}^{\tau}$ containing much less information than the trajectory itself. As indicated by Eq. (2), the KLD of work distributions should in principle be smaller than the KLD of trajectory distributions. On the contrary, the KLD is the same, indicating that all the irreversibility of the process is captured by the dissipative work [16].

C. Stationary trajectories

We now proceed to apply the above results to stationary trajectories. Consider a long process in which the system reaches a nonequilibrium stationary state (NESS) after a possible initial transient. In the NESS the external parameter is held fixed, $\lambda_t = \lambda$; the system is kept out of equilibrium due to the existence of baths at different temperatures (a possibility that is included in the hypothesis used in Ref. [2] to prove (3)) or different chemical potentials, external constant

forces, etc. In the steady state, since the control parameter remains fixed, the protocol and its time reversal are identical: $\lambda_t = \tilde{\lambda}_t = \lambda$ [13]. Therefore the probability distributions of the process and its time reversal are identical, $\tilde{\mathcal{P}} = \mathcal{P}$. In the long time limit, $\tau \rightarrow \infty$, we can neglect the contribution of the transient to the entropy production and rewrite (3) for the entropy production per unit of time \dot{S} in the NESS [28] as

$$\langle \dot{S} \rangle = \lim_{\tau \rightarrow \infty} \frac{k}{\tau} D[\mathcal{P}(\{z(t)\}_{t=0}^\tau) \| \mathcal{P}(\{\tilde{z}(\tau-t)\}_{t=0}^\tau)]. \quad (8)$$

A similar expression can be obtained from the Gallavotti-Cohen theorem [29,30], $\Delta S \simeq k \ln \frac{p_\tau(\Delta S)}{p_\tau(-\Delta S)}$, where $p_\tau(\Delta S)$ is the probability to observe an entropy production ΔS in the interval $[0, \tau]$. The Gallavotti-Cohen relationship, which is exact for $\tau \rightarrow \infty$, yields, after averaging,

$$\langle \dot{S} \rangle = \lim_{\tau \rightarrow \infty} \frac{k}{\tau} D[p_\tau(\Delta S) \| p_\tau(-\Delta S)]. \quad (9)$$

Consequently, although ΔS is another observable that is obtained as a function of the microstate of the system, the KLD calculated with ΔS yields the same value as the one calculated with full information of the system. Therefore entropy production captures all the information about the time irreversibility of the NESS.

When one does not observe the entire microscopic trajectory $\{z(t)\}_{t=0}^\tau$ in Eq. (8) but the trajectory followed by one or several observables of the system $x(t)$, the KLD only provides a lower bound to the entropy production [31]. Equations (7) and (9) indicate that the equality is recovered if the observables determine in a unique way the entropy production or the dissipated work.

In an experimental context, the observables are usually sampled at a finite frequency. The output is then a time series of data or discrete *trajectory*, $\mathbf{x} = (\hat{x}_1, \hat{x}_2, \dots, \hat{x}_n)$, where \hat{x}_i can be the value of a single or several observables of the system. In this case, we are interested in estimating the entropy production *per data* of the underlying physical process, which we denote by $\langle \dot{S} \rangle$ in the rest of the paper. Entropy production per data is related to the KLD rate *per data*, which we define below.

Given an infinitely long realization or time series sampled from a random process X_i ($i = 1, 2, \dots$), which can be multidimensional, we define by $p(x_1^m)$ the probability that a given string of m consecutive data is equal to $x_1^m = (x_1, x_2, \dots, x_m)$. We define the m th-order KLD for this random process X_i by the distinguishability between $p(x_1^m)$ and the probability $p(x_m^1)$ to observe the reverse sequence of data $x_m^1 = (x_m, x_{m-1}, \dots, x_1)$:

$$D_m^X = D[p(x_1^m) \| p(x_m^1)] = \sum_{x_1, \dots, x_m} p(x_1^m) \ln \frac{p(x_1^m)}{p(x_m^1)}. \quad (10)$$

The KLD rate for the process X_i is defined as the growth rate of D_m^X with the number of data:

$$d^X = \lim_{m \rightarrow \infty} \frac{D_m^X}{m}. \quad (11)$$

By virtue of (8) and (2), this quantity bounds from below the entropy production per data

$$\langle \dot{S} \rangle \geq k d^X, \quad (12)$$

where the bound is saturated if the random variable is the microstate of the system $X = \{\mathbf{q}, \mathbf{p}\}$ and the sampling rate is infinite [31] or X determines uniquely the entropy production in the process.

Equation (12) is our basic result. It reveals a striking connection between physics and the statistics of a time series. The left-hand side, $\langle \dot{S} \rangle$, is a purely physical quantity, whereas the right-hand side, d^X , is a statistical magnitude depending solely on the observed data, but not on the physical mechanism generating the data. Such a connection generalizes Landauer's principle relating entropy production and logical irreversibility in computing machines [1,32,33]. Equation (12) extends this principle and suggests that we can determine the average dissipation of an arbitrary NESS, even ignoring any physical detail of the system.

D. Markovian trajectories obeying local detailed balance

We first analyze how the bound (12) is expressed for Markovian time series that obey detailed balance by deriving analytical expressions for both entropy production and the KLD rate. If the random process X_i is Markovian, the probability distribution $p(x_1^m)$ factorizes $p(x_1^m) = p(x_1)p(x_2|x_1) \cdots p(x_m|x_{m-1})$, which also holds if we reverse the arguments, i.e., for $p(x_m^1)$. Substituting these expressions into Eq. (11), we get

$$d^X = \sum_{x_1, x_2} p(x_1, x_2) \ln \frac{p(x_2|x_1)}{p(x_1|x_2)} = D_2^X - D_1^X = D_2^X, \quad (13)$$

since $D_1^X = 0$ when comparing a trajectory and its reverse. Therefore, d^X only depends on transition probabilities if X is a random Markovian process.

We now relate d^X in Eq. (13) with the entropy production when the system reaches a NESS, because it is in contact with several thermal baths. In this situation, the local detailed balance condition is satisfied. We call $V(x_i)$ is the energy of the state x_i , and T_{x_1, x_2} is the temperature of the bath that activates the transitions $x_1 \rightarrow x_2$ and $x_2 \rightarrow x_1$. The local detailed balance condition reads in this case

$$\frac{p(x_2|x_1)}{p(x_1|x_2)} = \exp\left(\frac{V(x_1) - V(x_2)}{k T_{x_1, x_2}}\right). \quad (14)$$

Inserting (14) into (13),

$$\begin{aligned} d^X &= \sum_{x_1, x_2} p(x_1, x_2) \frac{V(x_1) - V(x_2)}{k T_{x_1, x_2}} \\ &= \sum_{x_1, x_2} p(x_1, x_2) \frac{Q_{x_1, x_2}}{k T_{x_1, x_2}} = \frac{\langle \dot{S} \rangle}{k}, \end{aligned} \quad (15)$$

where $Q_{x_1, x_2} = V(x_1) - V(x_2)$ is the heat dissipated to the corresponding thermal bath in the jump $x_1 \rightarrow x_2$, and \dot{S} is the total entropy production per data. Therefore, Eq. (12) is reproduced, with equality, in the case of a physical system obeying local detailed balance, if we have access to all the variables describing the system. The same conclusion is reached if we induce the NESS by means of nonconservative constant forces.

Equation (13) can be explored further by defining the *current* from the state x_1 to the state x_2 as the net probability flow from x_1 to x_2 , $J_{x_1 \rightarrow x_2} = p(x_1, x_2) - p(x_2, x_1)$. If the

system is not far from equilibrium the current tends to zero, and the condition $J_{x_1 \rightarrow x_2} \ll p(x_1, x_2)$ is satisfied, yielding

$$\frac{\langle \dot{S} \rangle}{k} = d^X = D_2^X \simeq \sum_{x_1, x_2} \frac{(J_{x_1 \rightarrow x_2})^2}{2p(x_1, x_2)}. \quad (16)$$

This expression is well known from linear irreversible thermodynamics [34], where entropy production is given by the product of a flow times a thermodynamic force that is proportional to the flow itself. Equation (16) implies that the time asymmetry of a Markovian process not far from equilibrium is revealed by the currents or probability flows that can be observed. In other words, a Markovian process without flows is time reversible. This is not the case for non-Markovian time series, where irreversibility can show up even in the absence of currents (see below and [9]).

III. KULLBACK-LEIBLER DIVERGENCE BETWEEN HIDDEN MARKOV CHAINS

In many experimental situations, a physical process is Markovian at a micro- or mesoscopic level of description, but the observed time series only contain a subset of the relevant observables, being non-Markovian in general. This is the case in biological systems, where one can only register the behavior of some mechanical and maybe a few chemical variables, while most of the relevant chemical variables cannot be monitored. These kinds of non-Markovian time series obtained from an underlying Markov process are called *hidden Markov chains* [35].

In this section we derive a semianalytical technique to calculate the KLD rate between hidden Markov chains. We focus on a simple case where the underlying Markov process is described by two observables X and Y ; however we only observe X whose evolution is described by a hidden Markov chain. The KLD rate for the observable X is

$$d^X = \lim_{m \rightarrow \infty} \frac{1}{m} \sum_{x_1^m} p(x_1^m) \ln \frac{\sum_{y_1^m} p(x_1^m, y_1^m)}{\sum_{y_1^m} p(x_1^m, y_1^m)}. \quad (17)$$

It is convenient to write d^X as a difference between two terms, $d^X = h_r^X - h^X$, where

$$h^X = - \lim_{m \rightarrow \infty} \frac{1}{m} \sum_{x_1^m} p(x_1^m) \ln \sum_{y_1^m} p(x_1^m, y_1^m) \quad (18)$$

is called the *Shannon entropy rate*, and

$$h_r^X = - \lim_{m \rightarrow \infty} \frac{1}{m} \sum_{x_1^m} p(x_1^m) \ln \sum_{y_1^m} p(x_1^m, y_1^m), \quad (19)$$

the *cross entropy rate*. Since the underlying process is Markovian, $p(x_1^m, y_1^m)$ factorizes and both Shannon and cross entropy can be expressed in terms of the trace of a product of random transition matrices \mathbf{T} [36,37]. These are square $M \times M$ random matrices, where M is the number of values that the variable y can take on, and their entries are given by

$$\mathbf{T}(x_1, x_2)_{y_1, y_2} = p(x_2, y_2 | x_1, y_1). \quad (20)$$

Note the different role played by each variable in this formalism: x_i are parameters defining the matrix (making \mathbf{T} a random matrix), whereas y_i are subindices of the matrix

elements. The Shannon and cross entropy can be expressed in terms of these matrices,

$$h^X = - \lim_{m \rightarrow \infty} \frac{1}{m} \left\langle \ln \text{Tr} \left[\prod_{i=1}^{m-1} \mathbf{T}(x_i, x_{i+1}) \right] \right\rangle, \quad (21)$$

$$h_r^X = - \lim_{m \rightarrow \infty} \frac{1}{m} \left\langle \ln \text{Tr} \left[\prod_{i=1}^{m-1} \mathbf{T}(x_{m-i+1}, x_{m-i}) \right] \right\rangle, \quad (22)$$

where $\langle \cdot \rangle$ denotes the average over the random process X_i , which are weighted by $p(x_1^m)$. For sufficiently large m , Eqs. (21) and (22) are self-averaging [37], meaning that we do not need to calculate the average but just compute the trace for a single stationary trajectory. For any sufficiently long time series $\mathbf{x} = (\hat{x}_1, \hat{x}_2, \dots, \hat{x}_n)$ with n large, the following expressions converge to $-h$ and $-h_r$ almost surely,

$$\hat{\lambda}^{\mathbf{x}} = \frac{1}{n} \ln \left\| \prod_{i=1}^{n-1} \mathbf{T}(\hat{x}_i, \hat{x}_{i+1}) \right\| \simeq -h^X, \quad (23)$$

$$\hat{\lambda}^{\tilde{\mathbf{x}}} = \frac{1}{n} \ln \left\| \prod_{i=1}^{n-1} \mathbf{T}(\hat{x}_{n-i+1}, \hat{x}_{n-i}) \right\| \simeq -h_r^X, \quad (24)$$

where $\| \cdot \|$ is any matrix norm that satisfies $\| \mathbf{A} \cdot \mathbf{B} \| \leq \| \mathbf{A} \| \| \mathbf{B} \|$ [37]. In particular, the trace satisfies this condition for positive matrices. In the context of random matrix theory, $\hat{\lambda}^{\mathbf{x}}$ and $\hat{\lambda}^{\tilde{\mathbf{x}}}$ are known as *maximum Lyapunov characteristic exponents* [38] and measure the asymptotic rate of growth of a random vector when being multiplied by a random sequence of matrices. In practice, we can estimate d^X semianalytically as

$$\hat{d}^{\mathbf{x}} = \hat{\lambda}^{\mathbf{x}} - \hat{\lambda}^{\tilde{\mathbf{x}}}. \quad (25)$$

Here $\hat{\lambda}^{\mathbf{x}}$ and $\hat{\lambda}^{\tilde{\mathbf{x}}}$ are estimated using (23) and (24) with a single time series \mathbf{x} of size n , following a technique introduced in Ref. [38]: We generate a random stationary time series $\mathbf{x} = \{\hat{x}_i^1\}$ and compute the matrices \mathbf{T} analytically; then a random unitary vector is multiplied by those matrices and normalized every l data, keeping track of the normalization factor; finally the product of these factors divided by n yields $\hat{\lambda}^{\mathbf{x}}$. For $\hat{\lambda}^{\tilde{\mathbf{x}}}$, the same procedure is repeated but using the reversed time series $\tilde{\mathbf{x}} = \{\hat{x}_n^1\}$. The technique is semianalytical since the transition probabilities are known analytically but a single random stationary time series \mathbf{x} is necessary to estimate d^X with the multiplication of n transition matrices that are chosen according to \mathbf{x} .

Let us recall that the estimator $\hat{d}^{\mathbf{x}}$ cannot be applied to empirical time series unless we know the Markov model behind the data. Consequently, it is not useful in practical situations. However, we will use it to check the performance of the estimators introduced in the following section, which only need a single stationary time series to estimate the KLD and do not assume any knowledge of the dynamics generating these data. On the other hand, one can also get analytical approximations of Eqs. (21) and (22) by using the replica trick, in a way analogous to that done in Ref. [39]. The calculation is cumbersome and is explained in the Appendix. Both the semianalytical and the replica calculations are used in Sec. V to check the accuracy of several empirical estimators of the KLD.

IV. ESTIMATING KLD RATES FROM SINGLE STATIONARY TRAJECTORIES

In previous sections, we calculated the KLD analytically (or semianalytically) for series where we know in advance the dynamics of the underlying physical process. We now investigate how the KLD rate can be estimated from a single empirical stationary trajectory, obtained from a discrete stochastic process whose dynamics is unknown. We call \hat{x}_i the value of the i th data of an empirical trajectory of n data, which is denoted by $\mathbf{x} = \{\hat{x}_i\}_{i=1}^n$. There are two types of estimators in the literature: *plug-in* estimators, based on empirical counting of sequences of data, and estimators based on compression algorithms. In this section, we introduce a refinement of the these two methods and analyze their performance for a specific example in Sec. V.

A. Plug-in estimators

The simplest approach to estimate the KLD rate is known as the *plug-in* method [24], which consists of an empirical estimation of the probabilities of sequences of m data, $p(x_1^m)$, appearing in Eq. (10). The probability of observing the sequence x_1^m , $p(x_1^m)$, is estimated empirically from simply counting the number of times that x_1^m appears in a single stationary trajectory $\mathbf{x} = (\hat{x}_1, \dots, \hat{x}_n)$ of size n . The empirical probability distribution is

$$\hat{p}^{\mathbf{x}}(x_1^m) = \frac{1}{n - (m - 1)} \sum_{p=1}^{n-(m-1)} \delta_{\hat{x}_p, x_1} \cdots \delta_{\hat{x}_{p+(m-1)}, x_m}. \quad (26)$$

Then an estimate of $D_m^{\mathbf{x}}$ is obtained by plugging the empirical probability distribution into Eq. (10):

$$\hat{D}_m^{\mathbf{x}} = D[\hat{p}^{\mathbf{x}}(x_1^m) \parallel \hat{p}^{\mathbf{x}}(x_1^1)] = \sum_{x_1, \dots, x_m} \hat{p}^{\mathbf{x}}(x_1^m) \ln \frac{\hat{p}^{\mathbf{x}}(x_1^m)}{\hat{p}^{\mathbf{x}}(x_1^1)}. \quad (27)$$

Note that the probabilities in Eq. (27) include the superscript \mathbf{x} to emphasize that they are obtained empirically from a single stationary time series \mathbf{x} and therefore depend on each particular realization. The simplest way to estimate $d^{\mathbf{x}}$ would be by taking $\frac{\hat{D}_m^{\mathbf{x}}}{m}$ for m as large as possible. However, this naive approach is not efficient. The empirical probability $\hat{p}^{\mathbf{x}}(x_1^m)$ —and therefore $\hat{D}_m^{\mathbf{x}}$ —is less accurate as m increases, because the number of possible substring x_1^m increases exponentially and the statistics shortly becomes poor. It is convenient to find alternative expressions with a fast convergence. It turns out that the slope of $\hat{D}_m^{\mathbf{x}}$ as a function of m ,

$$\hat{d}_m^{\mathbf{x}} = \hat{D}_m^{\mathbf{x}} - \hat{D}_{m-1}^{\mathbf{x}}, \quad (28)$$

also converges to the KLD rate but faster than $\frac{\hat{D}_m^{\mathbf{x}}}{m}$. Our plug-in estimator will be constructed as the limit

$$\hat{d}^{\mathbf{x}} = \lim_{m \rightarrow \infty} \hat{d}_m^{\mathbf{x}}. \quad (29)$$

For a Markovian time series, as shown in Eq. (13), the limit is reached for $m = 2$, and using distributions of three or more data we only get redundant information: $\hat{d}^{\mathbf{x}} = \hat{d}_2^{\mathbf{x}} = \hat{d}_m^{\mathbf{x}}$, for any $m > 2$. Therefore, $\hat{d}^{\mathbf{x}} = \hat{d}_2^{\mathbf{x}}$ is an excellent estimator of the KLD, $d^{\mathbf{x}}$. If \mathbf{x} is a k th-order Markov chain (i.e., it is Markovian when considering blocks of k data $\{\hat{x}_i^k\}$), then the limit is

reached for $m = k$, i.e., $\hat{d}^{\mathbf{x}} = \hat{d}_k^{\mathbf{x}} = \hat{d}_{k+1}^{\mathbf{x}} = \hat{d}_{k+2}^{\mathbf{x}} = \dots$ [23]. The convergence of (29) is then expected to be fast if a time series can be approximated by a k th-order Markov chain.

If the trajectory \mathbf{x} is sampled from a general non-Markovian process, one needs further information to extrapolate $\hat{d}_m^{\mathbf{x}}$ for $m \rightarrow \infty$, especially when only moderate values of m can be reached. In the examples discussed below, we have found that convergence is well described by the following ansatz, proposed by Schürmann and Grassberger [40] to estimate the Shannon entropy rate:

$$\hat{d}_m^{\mathbf{x}} \simeq \hat{d}_{\infty}^{\mathbf{x}} - c \frac{\ln m}{m^{\gamma}}. \quad (30)$$

Here c and γ are parameters that, together with $\hat{d}_{\infty}^{\mathbf{x}}$, can be obtained by fitting the empirical values of $\hat{d}_m^{\mathbf{x}}$ as a function of m . The fitting parameter $\hat{d}_{\infty}^{\mathbf{x}}$ gives an estimation of the limit (29).

This estimation method is efficient as long as there is sufficient statistics in the data, that is, if for every series x_1^m that occurs in the trajectory, its reverse x_m^1 is observed at least once. On the other hand, if we find empirically $\hat{p}^{\mathbf{x}}(x_1^m) \neq 0$ while $\hat{p}^{\mathbf{x}}(x_m^1) = 0$ for at least one case, the argument of the logarithm in Eq. (10) diverges, yielding $\hat{D}_m^{\mathbf{x}} = \infty$. We can avoid this divergence by restricting the sum in $\hat{D}_m^{\mathbf{x}}$ to sequences x_1^m whose reverse x_m^1 occurs in the time series:

$$\hat{D}_m^{\mathbf{x}} \rightarrow \hat{D}_m^{\mathbf{x}*} = \sum_{(x_1^m)^*} \hat{p}^{\mathbf{x}}(x_1^m) \ln \frac{\hat{p}^{\mathbf{x}}(x_1^m)}{\hat{p}^{\mathbf{x}}(x_m^1)}, \quad (31)$$

where $(x_1^m)^* = \{x_1^m \mid \hat{p}^{\mathbf{x}}(x_1^m) \neq 0 \text{ and } \hat{p}^{\mathbf{x}}(x_m^1) \neq 0\}$. With this restriction, a lower bound to $\hat{D}_m^{\mathbf{x}}$ is always obtained, $\hat{D}_m^{\mathbf{x}*} < \hat{D}_m^{\mathbf{x}}$.

A different strategy is to artificially bias the empirical probabilities such that all of them become positive. Instead of the observed empirical frequencies, we can use the following biased frequencies [41]:

$$\hat{p}^{\mathbf{x}}(x_1^m) = \frac{n^{\mathbf{x}}(x_1^m) + \gamma}{\sum_{x_1^m} [n^{\mathbf{x}}(x_1^m) + \gamma]}. \quad (32)$$

Here $n^{\mathbf{x}}(x_1^m)$ is the number of observations of x_1^m in \mathbf{x} and γ is the bias, which is a small number that prevents any of the probabilities from being zero, assigning a probability of order γ/n to sequences that are not observed. The denominator in Eq. (32) ensures normalization of $\hat{p}^{\mathbf{x}}(x_1^m)$.

B. Ziv-Merhav estimator

Ziv and Merhav introduced in Ref. [26] an estimator of the KLD rate between two probability distributions based on compression algorithms. It consists on slicing or *parsing* stationary discrete time series into smaller parts according to a specific algorithm. The slicing produces a sequence of numbers (often called a *dictionary*) that contains the same data as the original series, but it is divided into subsequences, called *phrases*. The algorithms that are used are called *compression algorithms* because the number of phrases in which a time series \mathbf{x} of n numbers is parsed into is smaller than n .

The estimator is defined in terms of two concepts which are now described, the compression length of a sequence and the

cross-parsing length between two different sequences. Given a series $\mathbf{x} = x_1^n$, its *compression length* $c(x_1^n)$ is defined as the number of distinct phrases in which it is parsed using the Lempel-Ziv (LZ) algorithm [42]. The LZ algorithm parses a series sequentially, such that each phrase that is added to the dictionary is the shortest distinct phrase that is not already in the dictionary. For example, let us consider the series $\mathbf{x} = x_1^{11} = (0, 1, 1, 1, 1, 0, 0, 1, 1, 0)$. The LZ sequential parsing for this example is as follows: First we store the first element of the sequence $x_1 = 0$ in the dictionary as it is empty, hence $\text{Dict} = \{0\}$. Then we read the next number, $x_2 = 1$, which is not already in the dictionary, so x_2 is added to the dictionary, $\text{Dict} = \{0|1\}$. The next number in x_1^{11} is $x_3 = 1$, which is already in the dictionary. Then we append to x_3 the next number of the sequence, $x_3^4 = (1, 1)$. This phrase is not in the dictionary and therefore it is parsed, $\text{Dict} = \{0|1|(1, 1)\}$. By doing this for all the series x_1^{11} , we obtain the following dictionary of phrases: $\text{Dict} = \{0|1|(1, 1)|(1, 0)|(0, 0)|(1, 1, 0)\}$. The compression length is the number of phrases that the dictionary contains once the series \mathbf{x} is completely parsed, $c(x_1^{11}) = 6$ in this example. The compression length of a stationary time series is related to its Shannon entropy rate [3] in the limit of infinitely long sequences:

$$\lim_{n \rightarrow \infty} \frac{c(x_1^n) \ln c(x_1^n)}{n} = h^X. \quad (33)$$

However, as $d^X = h_r^X - h^X$, we also require an estimator for h_r^X in order to determine d^X . This is given in terms of another quantity called the *cross-parsing length*. The cross parsing of a series x_1^n with respect to another sequence z_1^n is obtained by parsing x_1^n looking for the longest phrase that appears anywhere in z_1^n . As an example, let us consider the cross parsing of $\mathbf{x} = x_1^{11} = (0, 1, 1, 1, 1, 0, 0, 1, 1, 0)$ with respect to another sequence $\mathbf{z} = z_1^{11} = (1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 0)$. The first number in \mathbf{x} is $x_1 = 0$, which is in \mathbf{z} . Therefore we append to x_1 the next number in \mathbf{x} , $x_1^2 = (0, 1)$. This sequence is also somewhere in \mathbf{z} ; more precisely it is equal to z_3^4 , z_5^6 , and z_8^9 , so we append the next item in \mathbf{x} , $x_1^3 = (0, 1, 1)$. Again this sequence is somewhere in \mathbf{z} , $x_1^3 = z_8^{10}$, and it is added to the dictionary, $\text{Dict} = \{(0, 1, 1)\}$ because x_1^4 is not equal to any subsequence of z_1^{11} . We repeat this procedure again starting from x_4 and the resulting dictionary is $\text{Dict} = \{(0, 1, 1)|(1, 1, 0)|(0, 0, 1, 1, 0)\}$. The cross-parsing length is the number of parsed sequences, which in this example is equal to $c_r(x_1^{11}|z_1^{11}) = 3$. In Ref. [26] it is proved that the following quantity tends to the KLD rate between the probability distributions that generated the sequences $\mathbf{x} = x_1^n$ and $\mathbf{z} = z_1^n$, which we call p^X and q^Z , respectively:

$$\lim_{n \rightarrow \infty} \frac{1}{n} [c_r(x_1^n|z_1^n) \ln n - c(x_1^n) \ln c(x_1^n)] = d(p^X \| q^Z). \quad (34)$$

We can estimate d^X by using as inputs in the left-hand side of the above equation a stationary time series and its time reverse. The Ziv-Merhav estimator of d^X when using a time series \mathbf{x} of n data is introduced as follows:

$$\hat{d}_{ZM}^{\mathbf{x}} = \frac{1}{n} [c_r(x_1^n|x_n^1) \ln n - c(x_1^n) \ln c(x_1^n)], \quad (35)$$

which converges to d^X when $n \rightarrow \infty$, although the convergence is slow [26]. This estimator has been used as a measure

of distinguishability in several fields such as authorship attribution [22] or biometric identification [43].

When the KLD rate between the probability distributions under consideration is small ($d^X \ll 1$), the estimation given by Eq. (34) can be even negative [22]. The estimator gives negative values in some cases because it mixes two types of parsing: the sequential parsing of the trajectory and the cross parsing, which is not sequential. We propose the following correction, which helps to solve this issue and improves the performance of the estimator. We first evaluate (35) between different segments of the same trajectory. More precisely, we split \mathbf{x} into two equal parts and apply the original estimator (34)

$$\tilde{d}_{ZM}^{\mathbf{x}} = \frac{c_r(x_{n/2}^n|x_1^{n/2}) \ln \frac{n}{2} - c(x_{n/2}^n) \ln c(x_{n/2}^n)}{n/2}. \quad (36)$$

If the time series is stationary, the two fragments, $x_1^{n/2}$ and $x_{n/2}^n$, are equivalent and $\tilde{d}_{ZM}^{\mathbf{x}}$ should vanish. However it is usually negative for finite n and exhibits a slow convergence to zero for large n [22]. Then, we define our estimator as

$$\hat{d}_c^{\mathbf{x}} = \hat{d}_{ZM}^{\mathbf{x}} - \tilde{d}_{ZM}^{\mathbf{x}}, \quad (37)$$

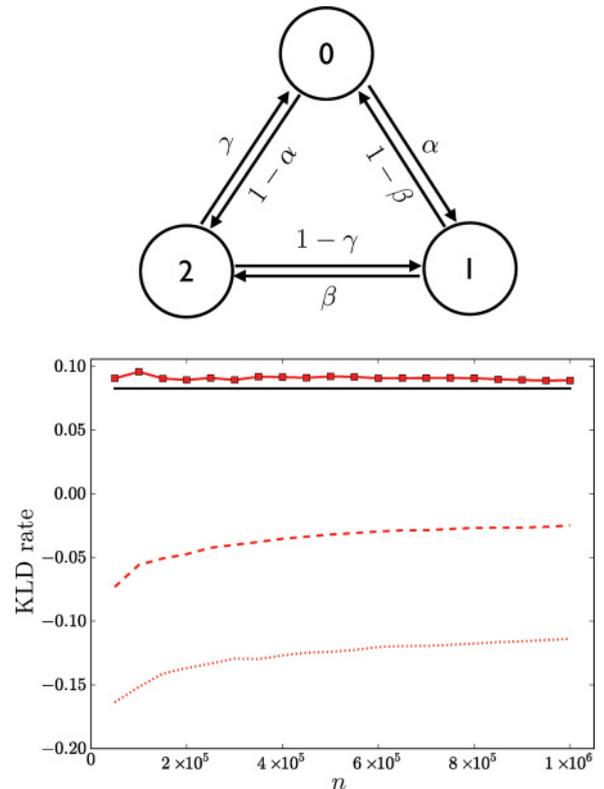


FIG. 1. (Color online) Sketch of the 3-state toy model used to check the accuracy of our compression estimator (37) and comparison between different compression estimators and the analytical value of d^X . The analytical value of d^X for a model with $\alpha = 0.5, \beta = 0.7, \gamma = 0.6$ ($d^X = d_2^X = 0.08278$) is indicated by the solid black line in the plot. We show the value of the compression estimators obtained from a single stationary time series x_1^n as a function of the length n : the Ziv-Merhav estimator $\hat{d}_{ZM}^{\mathbf{x}}$ (red dashed line), the bias $\tilde{d}_{ZM}^{\mathbf{x}}$ (red dotted line), and our estimator $\hat{d}_c^{\mathbf{x}}$ (red squares).

which still converges to d when $n \rightarrow \infty$ and yields much better results for finite n , as we show with a simple example.

We perform a first validation of this estimator using the three-state model illustrated in Fig. 1. Trajectories of the model are lists of numbers, 0, 1, or 2, representing the three states of the system. The dynamics is Markovian with transition probabilities given by $p_{0 \rightarrow 1} = 1 - p_{1 \rightarrow 0} = \alpha$, $p_{1 \rightarrow 2} = 1 - p_{2 \rightarrow 1} = \beta$, and $p_{2 \rightarrow 0} = 1 - p_{0 \rightarrow 2} = \gamma$. We call X_i the stochastic process describing the state of the system and \mathbf{x} a particular stationary time series, e.g., $\mathbf{x} = (0, 2, 1, 0, 1, 2, 1, 2, \dots)$. This time series is reversible only when the three transition probabilities satisfy the Kolmogorov condition [44], $\alpha\beta\gamma = (1 - \alpha)(1 - \beta)(1 - \gamma)$. In Fig. 1 (lower plot) we compare the value of different compression estimators with the analytical value of d^X as a function of the length of the empirical trajectory n . Since the trajectories described by the state of the system are Markovian, d^X only depends on transition probabilities: $d^X = d_2^X$. We see that the Ziv-Merhav estimator \hat{d}_{ZM}^x fails to estimate d^X accurately when it is small ($d^X \simeq 0.083$) and in some cases gives a negative value. The proposed estimator \hat{d}_c^x , on the other hand, is significantly closer to the analytical result, although it slightly overestimates its true value.

V. APPLICATION: THE DISCRETE FLASHING RATCHET

A. The model

We now apply the previous techniques to a specific example: a discrete flashing ratchet consisting of a Brownian particle moving on a one-dimensional lattice [45]. The particle is immersed in a thermal bath at temperature T and moves in a periodic, linear, asymmetric potential of height $2V$, which is switched on and off at a constant rate r (see Fig. 2). Trajectories are denoted by two random observables: the position of the particle X (0, 1, or 2) and the state of the potential Y (ON, $Y = 1$ or OFF, $Y = 0$).

The particle evolves in continuous time according to a master equation. The dynamics is described in terms of rates of spatial jumps and switching. For each possible transition except switches, i.e., $(x_1, y_1) \rightarrow (x_2, y_2)$ with $y_1 = y_2 = y$, we define a transition rate $k_{(x_1, y) \rightarrow (x_2, y)}$ obeying detailed balance,

$$k_{(x_1, y) \rightarrow (x_2, y)} = \exp\left[-\frac{V_y(x_2) - V_y(x_1)}{2kT}\right]. \quad (38)$$

When the potential is on ($y = 1$), the value of the potential energy $V_1(x)$ is given in Fig. 2. When the potential is off, $V_0(x) = 0$ for all x , and $k_{(x_1, 0) \rightarrow (x_2, 0)} = 1$ for $x_1 \neq x_2$. The switching rate does not depend on the position of the particle: $k_{(x, y_1) \rightarrow (x, y_2)} = r$ for any value of x and $y_1 \neq y_2$, and consequently violates detailed balance, driving the system out of equilibrium.

We simplify the analysis by mapping the dynamics onto a discrete-time process, a Markov chain. To this end, we record in a time series $(\mathbf{x}, \mathbf{y}) = \{x_i^n, y_i^n\}$ just a list of the visited states, discarding any information about the time where jumps and switches occur. The resulting Markov chain is defined by the transition probabilities

$$p[(x_2, y_2)|(x_1, y_1)] = \frac{k_{(x_1, y_1) \rightarrow (x_2, y_2)}}{\sum_{x_2, y_2} k_{(x_1, y_1) \rightarrow (x_2, y_2)}}. \quad (39)$$

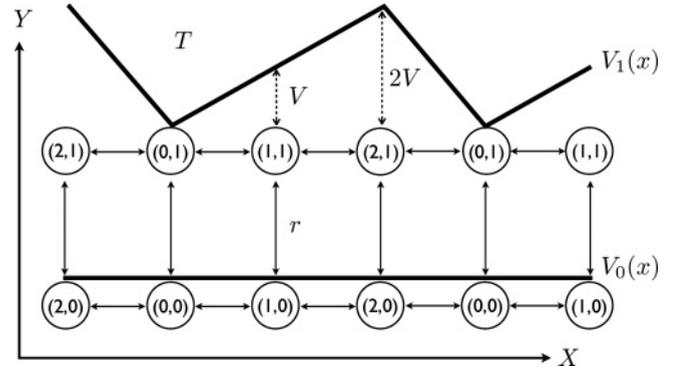


FIG. 2. Illustration of our discrete ratchet model. Particles are immersed in a thermal bath at temperature T and move in one dimension in an asymmetric linear potential $V_1(x)$ of height $2V$ with periodic boundary conditions. The potential is switched on and off at a rate r , where $V_0(x) = 0$ represents a flat potential, and the switching probability does not depend on the position of the particle. The state of the particle is represented by two random variables (X, Y) indicated in the figure, where $X = \{0, 1, 2\}$ stands for the position of the particle whereas $Y = \{0, 1\}$ for the state of the potential. Using this description, the system can be in six different states, $(0, 0)$, $(1, 0)$, $(2, 0)$, $(0, 1)$, $(1, 1)$, $(2, 1)$.

Since we discard any information about the transition times, we will focus in the rest of paper only on dissipation and KLD rates *per jump* or per data. For finite switching rate r , the ratchet rectifies the thermal fluctuations inducing a current to the left in Fig. 2 [34,45]. The system obeys a local detailed balance condition, as described in Sec. II D. The nonequilibrium nature of the switching can be interpreted in two alternative ways: One can imagine that it is activated by a thermal bath at infinite temperature or by an external agent [34]. In either of the two interpretations, switching does not induce any entropy production (the bath needs an infinite amount of energy to change its entropy and the external agent does not produce any entropy change). Therefore, entropy is only produced when heat is dissipated to the bath at temperature T , which only occurs when the potential is on. The average entropy production (or dissipation) per data in the time series is then [see (15)]

$$\langle \dot{S} \rangle = \sum_{y=0,1} \sum_{x_1, x_2=0,1,2} p[(x_1, y); (x_2, y)] \frac{V_y(x_1) - V_y(x_2)}{T}, \quad (40)$$

which is equal to the KLD rate when calculated for time series containing the information of both position and state of the system (which we call *full information*), $\langle \dot{S} \rangle = d^{X,Y} = d_2^{X,Y}$. We now analyze how can d be estimated using single stationary trajectories of this model, and how close this estimation is to the entropy production depending on the number of degrees of freedom of the system that are sampled in the time series.

B. Full information

First, we investigate the estimation of the KLD rate when using full information of the system (the position of the particle X and the state of the potential Y), and how close this KLD rate is to the actual entropy production of the process. In Fig. 3 we compare the actual dissipation and

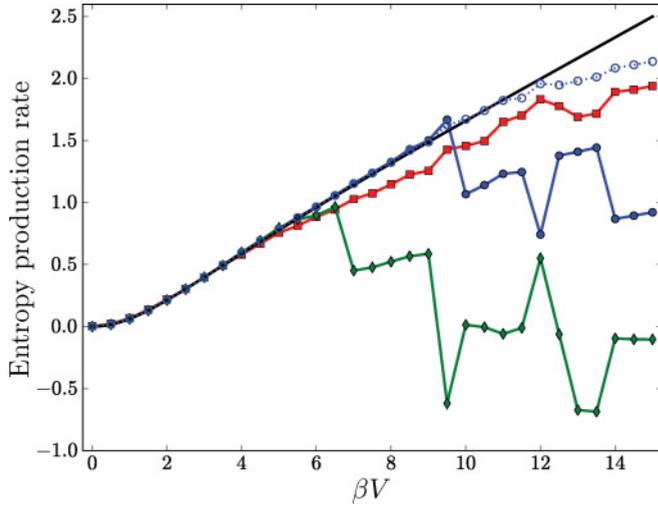


FIG. 3. (Color online) Analytical value of the average dissipation per data in units of kT (black line) as a function of βV in the flashing ratchet ($r = 1$) and different estimators of $d^{X,Y}$. For each value of βV , estimators are obtained from a single stationary time series of $n = 10^6$ data containing full information of the system (position, X , and state of the potential, Y). Plug-in estimators: $\hat{d}_2^{X,Y}$ (blue circles), $\hat{d}_3^{X,Y}$ (green diamonds), and $\hat{d}_2^{X,Y}$ using biased probabilities with $\gamma = 1$ (blue open circles). Compression estimator: $\hat{d}_c^{X,Y}$ (red squares).

several empirical estimations of $d^{X,Y}$ for different values of the height of the potential, V . For each value of V we simulate a single stationary time series of $n = 10^6$ data that contains full information, and calculate the plug-in estimators $\hat{d}_2^{X,Y}$, $\hat{d}_3^{X,Y}$, as well as the compression-based estimator $\hat{d}_c^{X,Y}$.

Since trajectories containing full information are Markovian, the plug-in estimator immediately converges to the dissipation $\hat{d}_2^{X,Y} = \hat{d}_3^{X,Y} = d^{X,Y} = \langle \dot{S} \rangle / k$ if there is enough statistics, which happens when V is below or of order kT . If $V \gg kT$, the observation of the uphill jumps such as $(0,1) \rightarrow (1,1)$, $(0,1) \rightarrow (2,1)$, or $(1,1) \rightarrow (2,1)$ is very unlikely in a single stationary trajectory. A time series of n data captures the statistics of jumps with probability well above $1/n$, which amounts to, say, energy jumps below $kT \ln n$ ($kT \ln 10^6 \approx 14kT$ for the trajectory used in the figures).

If, for instance, the transition $(0,1) \rightarrow (1,1)$ is missing in the trajectory, there is no way of estimating $p[(0,1);(1,1)]$ which contributes to two terms in $\hat{d}_2^{X,Y}$ [see Eq. (10) for $n = 2$]. One of these two terms accounts for jumps $(0,1) \rightarrow (1,1)$, which are very unlikely and their contribution to the total dissipation rate is negligible, and the other term accounts for jumps $(1,1) \rightarrow (0,1)$, whose probability is larger and therefore contributes more significantly to the entropy production.

In Fig. 3, $\hat{d}_2^{X,Y}$ (blue circles) and $\hat{d}_3^{X,Y}$ (green diamonds) have been calculated restricting the average to sequences (of two or three data respectively) whose reverse are also observed in the time series, as given by Eq. (31). The sudden drops in $\hat{d}_2^{X,Y}$ and $\hat{d}_3^{X,Y}$ are a consequence of lack of statistics in the trajectory. For the specific time series used in Fig. 3, the lack of statistics starts at $\beta V \simeq 10$ for $\hat{d}_2^{X,Y}$ and arises earlier for $\hat{d}_3^{X,Y}$ because

the three-data sampling space is bigger and it is easier that some transitions $(x_1, y_1) \rightarrow (x_2, y_2) \rightarrow (x_3, y_3)$ do not appear while their reverses do.

A more efficient way of dealing with the missing sequences is incorporating a small bias to the empirical probabilities, as described in Eq. (32). This is equivalent to assigning a probability of order $1/n$ to those transitions that are not observed in a time series of n data. Figure 3 shows $\hat{d}_2^{X,Y}$ with a bias $\gamma = 1$ (blue open circles), which is able to extend the accuracy of the estimation even when there is lack of statistics.

Although in the case of Markovian series with a finite number of states the most convenient strategy is to use the plug-in estimator, we include for comparison the compression estimator $\hat{d}_c^{X,Y}$ (red squares) which gives accurate values of the dissipation for weak potentials. Furthermore, the compression estimator is better than some plug-in estimators even for strong potentials, since it does not exhibit sudden jumps due to lack of statistics.

C. Partial information

We now analyze the performance of our estimators when there is no access to the full description of the system. As in Ref. [9], we assume that only the position of the ratchet X is observable. Accordingly, we simulate trajectories containing full information, and we remove the information of the state afterward, $(\mathbf{x}, \mathbf{y}) \rightarrow \mathbf{x}$. The resulting time series $\mathbf{x} = \{x_1^n\}$ is not Markovian and hence the limit (29) is not reached for small values of m . In this case, we proceed by obtaining \hat{d}_m^x for m as large as possible and fit the resulting values to the ansatz (30).

We have generated trajectories of size $n = 10^7$ for values of V that range from 0 to $2kT$. Once we remove the information of the state of the potential from these time series, we are able to estimate \hat{d}_m^x up to $m = 9$ with no lack of statistics. Figure 4 shows the plug-in estimators \hat{d}_m^x for $m = 2, 3, 5, 7, 9$ and the extrapolation \hat{d}_∞^x (orange pentagons connected by a dashed line to guide the eye) resulting from the fit to the ansatz (30). For each value of βV , we fit \hat{d}_m^x as a function of m for $m = 2, 3, \dots, 9$ to Eq. (30) using the curve-fitting tool available in MATLAB, which provides a robust least-squares fit with bisquare weights as described in Ref. [46]. The fit itself for a particular value of the potential, $\beta V = 1$, is shown in the inset of Fig. 4. Our ansatz reproduces the dependence of \hat{d}_m^x with m but the final estimator \hat{d}_∞^x still bounds significantly from below the actual dissipation (black solid line in Fig. 4). Nevertheless, plug-in estimators clearly distinguish between equilibrium and the NESS, even with partial information. In equilibrium ($V = 0$), the trajectories are reversible and all the estimators vanish, $\hat{d}_m^x = 0$ for $m = 2, \dots, 9$, whereas for the NESS ($V > 0$) they detect the irreversibility of the process yielding $\hat{d}_m^x > 0$ for all m . This is illustrated in Fig. 5, where we plot the dependence of the plug-in estimators with the size of the trajectory. For $\beta V = 0$, \hat{d}_2^x , \hat{d}_3^x , and \hat{d}_5^x tend to zero when increasing the number of data whereas they saturate to a positive value in the NESS ($\beta V = 1$).

There are two possible origins for the discrepancy between \hat{d}_∞^x and the dissipation: Either (i) our fit underestimates the actual KLD rate d^x of the trajectory; or (ii) the bound (12) is not tight. To address this question we need to calculate

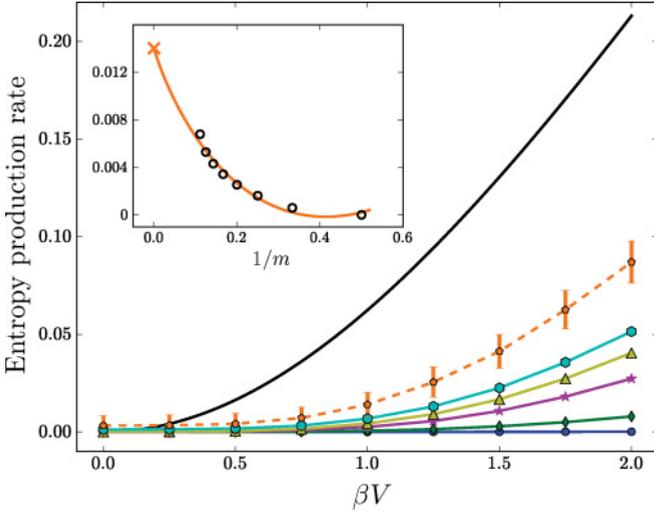


FIG. 4. (Color online) Average dissipation per data in units of kT (black line) and plug-in estimators of d^X using partial information given by the position (X) for a discrete flashing ratchet with $r = 1$. For each value of βV , we calculate estimators from a single stationary time series of $n = 10^7$ data containing partial information: \hat{d}_2^x (blue circles), \hat{d}_3^x (green diamonds), \hat{d}_5^x (purple stars), \hat{d}_7^x (yellow triangles), \hat{d}_9^x (cyan hexagons), and the result from the fit \hat{d}_∞^x (orange pentagons with error bars and connected by a dashed line). Inset: \hat{d}_m^x as a function of $1/m$ for $m = 1, \dots, 9$ for $\beta V = 1$ (open black circles) and the fit to the ansatz (orange line). The y intercept of the fit is indicated by an orange cross and it is equal to \hat{d}_∞^x .

the actual value of d^X . Since the position of the ratchet \mathbf{x} is a hidden Markov chain, we can calculate its KLD rate d^X semianalytically, using the Lyapunov exponents (23), (24) introduced in Sec. III.

In Fig. 6 we show the value of the semianalytical calculation of d^X using the norm of transition matrices, Eq. (25), which is not significantly different from the empirical estimation \hat{d}_∞^x . We therefore conclude that \hat{d}_∞^x is a good estimation of d^X , but still d^X only yields a lower bound to dissipation whose accuracy is in principle hard to determine. This is an expected

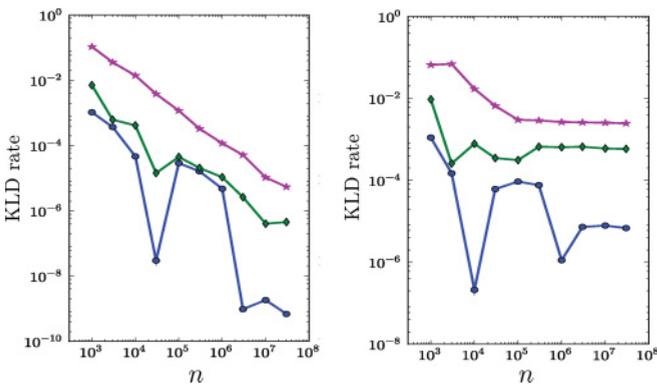


FIG. 5. (Color online) Scaling of plug-in estimators of d^X , \hat{d}_m^x , with the size of the time series n , for a flashing ratchet ($r = 1$), for $\beta V = 0$ (left) and $\beta V = 1$ (right): \hat{d}_2^x (blue circles), \hat{d}_3^x (green diamonds), and \hat{d}_5^x (purple stars). We simulate a single stationary trajectory \mathbf{x} of 10^7 data and calculate the estimators for subsequent ones containing the first n data of \mathbf{x} .

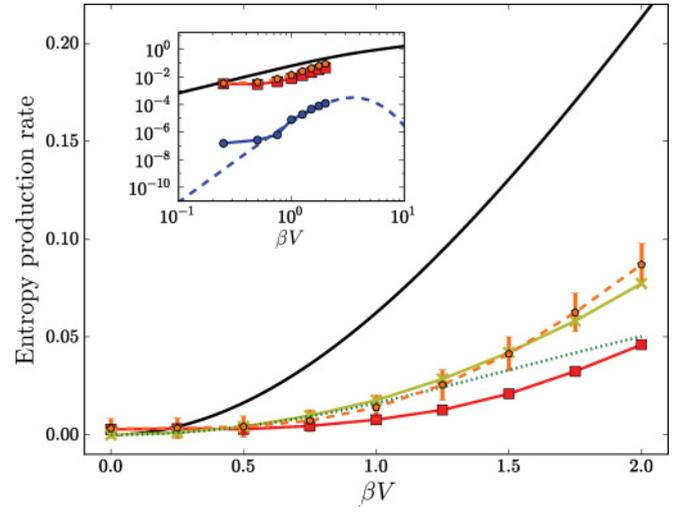


FIG. 6. (Color online) Average dissipation per data in units of kT (black line) and different estimators of d^X for a flashing ratchet described with partial information ($r = 1$, $n = 10^7$ data) as a function of βV : \hat{d}_∞^x (orange dashed pentagons), \hat{d}_c^x (red squares), replica estimation of d^X (green dotted line), and semianalytical value of d^X (yellow crosses). Inset: Dependence of the average dissipation (black line), \hat{d}_2^x (analytical values in blue dashed line), \hat{d}_c^x , and \hat{d}_∞^x on βV in the vicinity of $\beta V = 0$.

result, since the position of a particle in a flashing ratchet does not obey the Gallavotti-Cohen theorem [47].

Summarizing, although \hat{d}_∞^x turns out to be a good estimator of d^X , using only information of the position we only get a lower bound to the dissipation. We also show in Fig. 6 the value of \hat{d}_c^x , which is well below the plug-in estimator \hat{d}_∞^x . The compression estimator \hat{d}_c^x lies between \hat{d}_2^x and \hat{d}_5^x (not shown in the plot), indicating that it is only able to capture correlations up to size 8. For completeness, we include the calculation of d^X based on the replica trick (see the Appendix). It yields a tight bound for $V < kT$, but departs from d^X for larger values of V . This deviation is caused by the estimation of the limits in Eqs. (A10) and (A13), where we take $\alpha \rightarrow 0$ when α is defined only for integer values, one of the standard drawbacks of the replica trick [38].

Although our estimators give low values of the dissipation when using partial information, they still capture the asymptotic behavior for V small. Entropy production decreases as V^2 when $V \rightarrow 0$, so do plug-in estimators $\hat{d}_3^x, \dots, \hat{d}_9^x, \hat{d}_\infty^x$, and the compression estimator \hat{d}_c^x . Some of them are plotted in the inset of Fig. 6 (inset). On the other hand, $\hat{d}_2^x \propto V^6$, since the current is $J \propto V^3$ in this case [see Eq. (16)]. Recall that calculating \hat{d}_2^x is equivalent to estimating the entropy production using currents and standard linear irreversible thermodynamics, as shown in Eq. (16). It is then remarkable that the estimators involving the statistics of three or more data are able to reproduce qualitatively the behavior of the dissipation in cases where linear thermodynamics fails.

The improvement observed when using the plug-in estimators of higher order than \hat{d}_2^x is more dramatic in a NESS which does not exhibit observable currents in X . In this case $\hat{d}_2^x = 0$ but using higher order statistics we can still detect the time irreversibility of the trajectory [9]. This happens for example

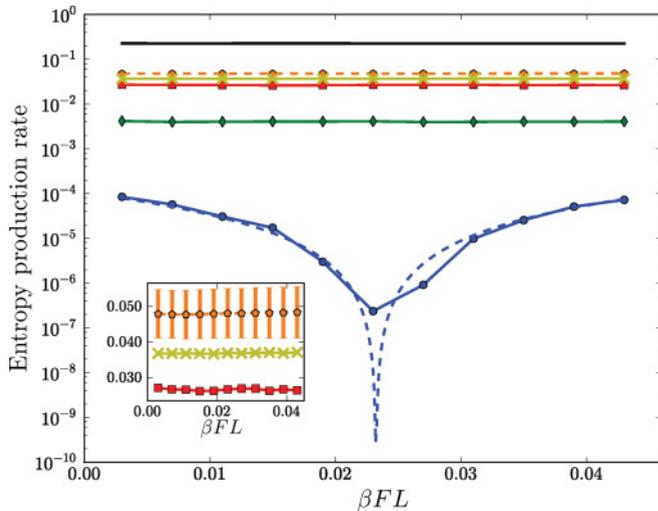


FIG. 7. (Color online) Average dissipation per data (in units of kT) in the flashing ratchet (with $r = 2$ and $\beta V = 2$) and different estimations of d^X obtained from a single time series of $n = 10^7$ data containing partial information (position) as a function of the external force F : analytical value of the average dissipation (black line), \hat{d}_2^X (blue circles, analytical values in blue dashed line), \hat{d}_3^X (green diamonds), \hat{d}_c^X (red squares), semianalytical calculation of d^X (yellow crosses), and \hat{d}_∞^X (orange hexagons). The minimum in \hat{d}_2^X corresponds to the stalling force. Inset: \hat{d}_c^X , semianalytical value of d^X and \hat{d}_∞^X as a function of the external force.

if we add to the flashing ratchet an external force F opposite to the current, i.e., pointing in the positive x direction. The force modifies the energy landscape and consequently the spatial transition rates $k_{(x_1,y) \rightarrow (x_2,y)}$ by a factor $\exp[\beta FL_{(x_1,y);(x_2,y)}/2]$, $L_{(x_1,y);(x_2,y)}$ being the spatial distance that separates the two points (x_1, y) and (x_2, y) . Here $L_{(x_1,y);(x_2,y)}$ is defined positive if the jump $(x_1, y) \rightarrow (x_2, y)$ points in the same direction as the force (i.e., to the right), and negative otherwise. At the *stalling force* F_{stall} , the current is canceled by the force and the system does not move on average when it is described only by X , but still dissipates energy. If we only have access to the information of the position, the system looks like it is in equilibrium: The spatial current vanishes, and so does \hat{d}_2^X , as shown in Fig. 7. However, there is a finite dissipation (black line in the figure) and the corresponding irreversibility is captured by the statistics of substrings of length 3 or more. Although d^X is below the real dissipation by an order of magnitude (see the semianalytical value of d^X ; yellow crosses in Fig. 7), it does not exhibit any sensible change at stall force. Finally, both \hat{d}_∞^X and \hat{d}_c^X provide estimates of d^X which are correct within one order of magnitude (see the inset of Fig. 7).

VI. CONCLUSIONS

We have shown that it is possible to estimate the entropy production rate by analyzing statistical properties of a time series observed in a NESS. The Kullback-Leibler divergence (KLD) per data between the time series and its time reversal is a lower bound to the entropy production rate.

We have introduced two estimators of this KLD rate, one based on empirical frequencies and another on compression algorithms, and we have checked their performance in a

specific example: a discrete flashing ratchet. We show that the KLD is a powerful tool to identify nonequilibrium states and to estimate the entropy production of a process, if this entropy production is of order of the Boltzmann constant. We have also shown that the bound given by the KLD can detect a nonzero dissipation even when the data do not exhibit any measurable flows.

Let us summarize our results by presenting a “recipe” to estimate the KLD from an experimental time series recorded from a discrete system in a NESS. If the number of possible states of the system is small enough, the best approach is to calculate the plug-in estimators \hat{d}_m^X (28) and then check the convergence when m increases. The possible lack of statistics can be circumvented using a small artificial bias, as discussed in Sec. IV A. If \hat{d}_m^X saturates for some value m^* , then the time series is an m^* -th-order Markov process and $\hat{d}^X = \hat{d}_{m^*}^X$. Otherwise, we can use the ansatz (30) and obtain \hat{d}_∞^X which is a good estimate of the KLD rate.

A second and complementary approach is the use of the compression estimator introduced in Sec. IV B. The estimator yields correct results in the examples that we have analyzed, but there is no clue about the corresponding error. Nevertheless, the compression estimator could be the only possible approach if the number of states of the time series is large. In this case, the calculation of empirical probability distributions $\hat{p}(x_1^m)$ would be unfeasible even for short substrings.

Another possible strategy for systems with many states (or described by real-valued observables) is to consider time-asymmetric functionals of the data, which reduce the number of observables, and hence the number of states, but keep information about the irreversibility of the series. In any case, the estimation of KLD and the extension of our results to processes described by continuous data is an open problem, which will be relevant in many practical situations, especially to analyze data coming from biological systems.

Finally, let us mention that, as in the case of Landauer’s principle, the KLD could also be used to ascertain the minimal entropy production associated with a specific behavior, such as spatiotemporal patterns and excitable systems. This in turn may influence the design of optimal devices with functionalities given by these behaviors.

ACKNOWLEDGMENTS

We acknowledge fruitful discussions with J. M. Horowitz and financial support from Grants MOSAICO (Spanish Government) and MODELICO (Comunidad Autónoma de Madrid, Spain). E.R. acknowledges funding from *Becas de la Caixa para estudios de Máster en España* and *Beca FPU* (Spanish Government).

APPENDIX: CALCULATION OF THE KLD RATE FOR HIDDEN MARKOV CHAINS USING REPLICA TRICK

The semianalytical calculation of the KLD rate for a specific case of hidden Markov chains was discussed in Sec. III. We now introduce another technique to calculate Eq. (17) using a mathematical technique called the *replica trick*. To this end, we first consider the expression of d^X in terms of Shannon

and cross entropy rates, $d^X = h_r^X - h^X$. We define the matrix resulting from the multiplication of m transition matrices [defined in Eq. (20)] chosen according to x_1^m by

$$\mathbf{T}(x_1^m) = \prod_{i=1}^{m-1} \mathbf{T}(x_i, x_{i+1}). \quad (\text{A1})$$

The Shannon entropy rate h^X can be rewritten by substituting (A1) into Eq. (21),

$$h^X = - \lim_{m \rightarrow \infty} \frac{1}{m} \langle \ln \text{Tr} \mathbf{T}(x_1^m) \rangle. \quad (\text{A2})$$

The analytical calculation of the average $\langle \ln \text{Tr} \mathbf{T}(x_1^m) \rangle$ is cumbersome and it can only be done semianalytically, as we explained in Sec. III. However, we can express this average in terms of $\langle \text{Tr} \mathbf{T}(x_1^m) \rangle$, which can be calculated analytically. The mathematical technique to do this is called the replica trick and it was introduced to calculate free energies in spin glasses [48]. For our specific example, the trick is given by the following expression:

$$\langle \ln \text{Tr} \mathbf{T}(x_1^m) \rangle = \lim_{\alpha \rightarrow 0} \frac{d}{d\alpha} \ln \langle [\text{Tr} \mathbf{T}(x_1^m)]^\alpha \rangle. \quad (\text{A3})$$

Reference [39] shows how to apply this technique when $\mathbf{T}(x_1^m)$ is equal to a product of random matrices which are chosen following a Markovian process. In our case, an underlying Markovian process defined by two random variables, X and Y , defines the order of the matrices that are multiplied in $\mathbf{T}(x_1^m)$. We now apply the technique described in Ref. [39] to calculate h^X . If we define the *generalized Lyapunov exponent* of degree α [38] as

$$L_\alpha^X = \lim_{m \rightarrow \infty} \frac{1}{m} \ln \langle [\text{Tr} \mathbf{T}(x_1^m)]^\alpha \rangle, \quad (\text{A4})$$

and we take into account the replica trick (A3), the Shannon entropy rate (A2) is given by

$$h^X = - \lim_{\alpha \rightarrow 0} \frac{d}{d\alpha} L_\alpha^X. \quad (\text{A5})$$

Now we consider the following property: Given a matrix \mathbf{A} and a positive integer α , $(\text{Tr} \mathbf{A})^\alpha = \text{Tr}(\mathbf{A}^{\otimes \alpha})$, where $\mathbf{A}^{\otimes \alpha} = \underbrace{\mathbf{A} \otimes \mathbf{A} \otimes \dots \otimes \mathbf{A}}_{\alpha \text{ times}}$. Using this property, the average in Eq. (A4) reads

$$\langle [\text{Tr} \mathbf{T}(x_1^m)]^\alpha \rangle = \langle \text{Tr}[\mathbf{T}(x_1^m)^{\otimes \alpha}] \rangle = \text{Tr}(\mathbf{T}(x_1^m)^{\otimes \alpha}). \quad (\text{A6})$$

Since the tensor power of a product of matrices factorizes, $(\mathbf{ABC})^{\otimes \alpha} = \mathbf{A}^{\otimes \alpha} \mathbf{B}^{\otimes \alpha} \mathbf{C}^{\otimes \alpha}$, Eq. (A6) can be

rewritten

$$\langle [\text{Tr} \mathbf{T}(x_1^m)]^\alpha \rangle = \text{Tr} \sum_{x_1^m, y_1^m} \prod_{i=1}^{m-1} \mathbf{T}(x_i, x_{i+1})_{y_i, y_{i+1}} \mathbf{T}(x_i, x_{i+1})^{\otimes \alpha}. \quad (\text{A7})$$

We now define a block matrix $\mathcal{T}(\alpha)$, where each block is a transition matrix $\mathbf{T}(x_1, x_2)^{\otimes \alpha+1}$. The matrix elements of $\mathcal{T}(\alpha)$ are therefore

$$\mathcal{T}(\alpha)_{x_1, y_1, x_2, y_2} = [\mathbf{T}(x_1, x_2)^{\otimes \alpha+1}]_{y_1, y_2}. \quad (\text{A8})$$

Using (A7) and (A8) in Eq. (A4), we see that L_α^X is dominated by the largest eigenvalue of $\mathcal{T}(\alpha)$ which we call $\tau(\alpha)$,

$$L_\alpha^X = \lim_{m \rightarrow \infty} \frac{1}{m} \ln \text{Tr}[\mathcal{T}(\alpha)^{m-1}] = \ln \tau(\alpha), \quad (\text{A9})$$

yielding

$$h^X = - \lim_{\alpha \rightarrow 0} \frac{d}{d\alpha} \ln \tau(\alpha). \quad (\text{A10})$$

The above limit cannot be calculated analytically because the tensor powers in $\mathcal{T}(\alpha)$ are only defined for integer values of α . Therefore we approximate the limit $\alpha \rightarrow 0$ by an estimation of the slope of L_α^X as a function of α close to $\alpha = 0$, which is given by [38]

$$\hat{h}^X = 2L_1^X - \frac{L_2^X}{2} = 2 \ln \tau(1) - \frac{\ln \tau(2)}{2}. \quad (\text{A11})$$

We obtain an equivalent result for h_r^X by replacing $\mathbf{T}(x_1^m)$ in Eq. (A2) by the product of transition matrices but ordered according to the time-reversed series $x_m^1, \mathbf{T}(x_m^1)$. Defining the matrix

$$\mathcal{T}_r(\alpha)_{x_1, y_1, x_2, y_2} = [\mathbf{T}(x_2, x_1)^T \otimes \mathbf{T}(x_1, x_2)^{\otimes \alpha}]_{y_1, y_2}, \quad (\text{A12})$$

and $\tau_r(\alpha)$ being the largest eigenvalue of $\mathcal{T}_r(\alpha)$, we get

$$h_r^X = - \lim_{\alpha \rightarrow 0} \frac{d}{d\alpha} \ln \tau_r(\alpha). \quad (\text{A13})$$

In practice, we also need to approximate the limit $\alpha \rightarrow 0$ in the above equation using Eq. (A11) but replacing τ by τ_r ,

$$\hat{h}_r^X = 2 \ln \tau_r(1) - \frac{\ln \tau_r(2)}{2}. \quad (\text{A14})$$

Finally, the estimation of d^X for this kind of series using the replica trick, which is shown in Fig. 6 (green dotted line), is obtained with the difference between Eqs. (A14) and (A11),

$$\hat{d}^X = \hat{h}_r^X - \hat{h}^X = 2 \ln \frac{\tau_r(1)}{\tau(1)} + \frac{1}{2} \ln \frac{\tau(2)}{\tau_r(2)}. \quad (\text{A15})$$

[1] R. Kawai, J. M. R. Parrondo, and C. Vanden Broeck, *Phys. Rev. Lett.* **98**, 080602 (2007).
 [2] J. M. R. Parrondo, C. V. den Broeck, and R. Kawai, *New J. Phys.* **11**, 073008 (2009).
 [3] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. (Wiley, Hoboken, NJ, 2006).
 [4] G. E. Crooks and D. A. Sivak, *J. Stat. Mech.: Theory Exp.* (2011) P06003.

[5] C. Maes and K. Netocny, *J. Stat. Phys.* **110**, 269 (2003).
 [6] C. Jarzynski, *Phys. Rev. E* **73**, 046105 (2006).
 [7] P. Gaspard, *J. Stat. Phys.* **117**, 599 (2004).
 [8] P. Zolfaghari, S. Zare, and B. Mirza, *Phys. Rev. E* **82**, 052104 (2010).
 [9] E. Roldán and J. M. R. Parrondo, *Phys. Rev. Lett.* **105**, 150607 (2010).
 [10] J. L. Lebowitz and H. Spohn, *J. Stat. Phys.* **95**, 333 (1999).

- [11] O. Mazonka and C. Jarzynski, e-print [arXiv:cond-mat/9912121](https://arxiv.org/abs/cond-mat/9912121).
- [12] R. van Zon, S. Ciliberto, and E. G. D. Cohen, *Phys. Rev. Lett.* **92**, 130601 (2004).
- [13] U. Seifert, *Phys. Rev. Lett.* **95**, 040602 (2005).
- [14] D. Andrieux, P. Gaspard, S. Ciliberto, N. Garnier, S. Joubaud, and A. Petrosyan, *Phys. Rev. Lett.* **98**, 150601 (2007).
- [15] J. Horowitz and C. Jarzynski, *Phys. Rev. E* **79**, 021106 (2009).
- [16] A. Gomez-Marin, J. M. R. Parrondo, and C. V. den Broeck, *Europhys. Lett.* **82**, 50002 (2008).
- [17] J. Kurchan, *J. Phys. A* **31**, 3719 (1998).
- [18] N. Garnier and S. Ciliberto, *Phys. Rev. E* **71**, 060101 (2005).
- [19] P. Martin, A. J. Hudspeth, and F. Jülicher, *Proc. Natl. Acad. Sci. USA* **98**, 14380 (2001).
- [20] C. P. Amann, T. Schmiedl, and U. Seifert, *J. Chem. Phys.* **132**, 041102 (2010).
- [21] M. B. Kennel, *Phys. Rev. E* **69**, 056208 (2004).
- [22] D. P. Coutinho and M. A. Figueiredo, *Pattern Recognition and Image Analysis*, Vol. 3523 of Lecture Notes in Computer Science (Springer, Berlin, 2005).
- [23] Z. Rached, F. Alajaji, and L. L. Campbell, *IEEE Trans. Inf. Theory* **50**, 917 (2004).
- [24] Q. Wang, S. Kulkarni, and S. Verdu, *IEEE Trans. Inf. Theory* **51**, 3064 (2005).
- [25] M. Budka, B. Gabrys, and K. Musial, *Entropy* **13**, 1229 (2011).
- [26] J. Ziv and N. Merhav, *IEEE Trans. Inf. Theory* **39**, 1270 (1993).
- [27] G. E. Crooks, *Phys. Rev. E* **60**, 2721 (1999).
- [28] C. Maes, *Sem. Poincaré* **2**, 29 (2003).
- [29] E. Cohen and G. Gallavotti, *J. Stat. Phys.* **96**, 1343 (1999).
- [30] B. Cleuren, K. Willaert, A. Engel, and C. Van den Broeck, *Phys. Rev. E* **77**, 022103 (2008).
- [31] A. Gomez-Marin, J. M. R. Parrondo, and C. Van den Broeck, *Phys. Rev. E* **78**, 011107 (2008).
- [32] R. Landauer, *IBM J. Res. Dev.* **44**, 261 (2000).
- [33] D. Andrieux and P. Gaspard, *Proc. Natl. Acad. Sci. USA* **105**, 9516 (2008).
- [34] J. M. R. Parrondo and B. J. de Cisneros, *Appl. Phys. A* **75**, 179 (2002).
- [35] L. Rabiner and B. Juang, *IEEE ASSP Magazine* **3**, 4 (1986).
- [36] P. Jacquet, G. Seroussi, and W. Szpankowski, *Theor. Comput. Sci.* **395**, 203 (2008).
- [37] P. G. T. Holliday and A. Goldsmith, *IEEE Trans. Inf. Theory* **52**, 3509 (2006).
- [38] A. Crisanti, G. Paladin, and A. Vulpiani, *Products of Random Matrices in Statistical Physics* (Springer Series in Solid State Sciences, 1993).
- [39] M. J. de Oliveira and A. Petri, *Phys. Rev. E* **53**, 2960 (1996).
- [40] T. Schurmann and P. Grassberger, *Chaos* **6**, 414 (1996).
- [41] H. Cai, S. Kulkarni, and S. Verdu, *IEEE Trans. Inf. Theory* **52**, 3456 (2006).
- [42] J. Ziv and A. Lempel, *IEEE Trans. Inf. Theory* **24**, 530 (1978).
- [43] D. P. Coutinho, A. L. Fred, and M. A. Figueiredo, *20th International Conference on Pattern Recognition*, 3858 (2010).
- [44] J. F. C. Kingman, *J. Appl. Prob.* **6**, 1 (1969).
- [45] A. Ajdari and J. Prost, *C. R. Acad. Sci. Paris Ser. II* **315**, 1635 (1992).
- [46] MATLAB R2011b documentation (curve fitting toolbox), <http://www.mathworks.es/help/toolbox/curvefit/>.
- [47] D. Lacoste and K. Mallick, *Phys. Rev. E* **80**, 021923 (2009).
- [48] J. van Hemmen and R. Palmer, *J. Phys. A* **12**, 563 (1979).