

Fluctuations When Driving Between Nonequilibrium Steady States

Paul M. Riechers
James P. Crutchfield

SFI WORKING PAPER: 2016-10-023

SFI Working Papers contain accounts of scientific work of the author(s) and do not necessarily represent the views of the Santa Fe Institute. We accept papers intended for publication in peer-reviewed journals or proceedings volumes, but not papers that have already appeared in print. Except for papers by our external faculty, papers must be based on work done at SFI, inspired by an invited visit to or collaboration at SFI, or funded by an SFI grant.

©NOTICE: This working paper is included by permission of the contributing author(s) as a means to ensure timely distribution of the scholarly and technical work on a non-commercial basis. Copyright and all rights therein are maintained by the author(s). It is understood that all persons copying this information will adhere to the terms and constraints invoked by each author's copyright. These works may be reposted only with the explicit permission of the copyright holder.

www.santafe.edu



SANTA FE INSTITUTE

Fluctuations When Driving Between Nonequilibrium Steady States

Paul M. Riechers* and James P. Crutchfield†

*Complexity Sciences Center, Department of Physics
University of California at Davis, One Shields Avenue, Davis, CA 95616*

(Dated: October 29, 2016)

Maintained by environmental fluxes, biological systems are thermodynamic processes that operate far from equilibrium without detailed-balance dynamics. Yet, they often exhibit well defined nonequilibrium steady states (NESSs). More importantly, critical thermodynamic functionality arises directly from transitions among their NESSs, driven by environmental switching. Here, we identify constraints on excess thermodynamic quantities that ride above the NESS housekeeping background. We do this by extending the Crooks fluctuation theorem to transitions among NESSs, without invoking an unphysical dual dynamics. This and corresponding integral fluctuation theorems determine how much work must be expended when controlling systems maintained far from equilibrium. This generalizes feedback control theory, showing that Maxwellian Demons can leverage mesoscopic-state information to take advantage of the excess energetics in NESS transitions. Altogether, these point to universal thermodynamic laws that are immediately applicable to the accessible degrees of freedom within the effective dynamic at any emergent level of hierarchical organization. By way of illustration, this readily allows analyzing a voltage-gated sodium ion channel whose molecular conformational dynamics play a critical functional role in propagating action potentials in mammalian neuronal membranes.

PACS numbers: 05.70.Ln 05.20.-y 05.45.-a 02.50.Ey

Keywords: stochastic thermodynamics, fluctuation theorem, nonequilibrium, neuronal ion channel

I. INTRODUCTION

The sun shines; ATP is abundant; power is supplied. These are the generous settings in which we find many complex biological systems, buoyed steadily out of equilibrium by energy fluxes in their environment. The resulting steady-state dynamics exhibit various types of directionality, including periodic oscillations and macroscopic thermodynamic functionality. These behaviors contrast rather sharply with the deathly isotropy of equilibrium detailed-balance dynamics—where fluxes are absent and state transition rates depend only on relative asymptotic state-occupation probabilities.

Detailed balance and its implied dynamical reversibility, though, are common tenets of equilibrium thermodynamics and statistical mechanics. They are technically necessary when applying much of the associated theory relevant to equilibrium and reversible (e.g., quasistatic) transitions between equilibrium macrostates [1]. Detailed balance is even assumed by several modern theorems that influence our understanding of the structure of fluctuations and limitations on work performed far from equilibrium [2, 3]. The natural world, though, is replete with systems that violate detailed balance, such as small biological molecules constantly driven out of equilibrium through interactions with their biochemical environment [4, 5].

Far from happenstance and disruption, the probability currents through the effective state-space of these nonequilibrium systems enable crucial thermodynamic functionality [6, 7]. Even rare fluctuations play an important functional role [8–10]. While constant environmental pressure can drive a system into a nonequilibrium steady state (NESS), complex biological systems are often driven farther—far from even any NESS. Moreover, such system–environment dynamics involve feedback between system and environment states. Although many believe these facilitate the necessary complex processes that sustain life, their very nature seems to preclude most, if not all, hope of a universal theoretical framework for quantitative predictions. To ameliorate the roadblock, we present a consistent thermodynamics that is not only descriptive, but constructive, tractable, and predictive, even when irreversible dynamics transition between NESSs.

Beyond laying out the structure of fluctuations during NESS transitions, this thermodynamics sets the stage to understand how one level of organization gives way to another. In particular, using it a sequel renormalizes nonequilibrium housekeeping background to show how to maintain a hierarchy of steady-state dynamics. Said simply, at each level of hierarchical organization, controllable degrees of freedom are subject to universal thermodynamic laws that tie their fluctuations and functionality to dissipation at lower levels.

* pmriechers@ucdavis.edu

† chaos@ucdavis.edu

A. Results

Nonequilibrium thermodynamics progressed markedly over the last two decades on at least two fronts. First, by taking the ‘dynamics’ in ‘thermodynamics’ seriously, fluctuation theorems (FTs) transformed previous inequalities, such as the classical Second Law of Thermodynamics, into subsuming equalities that exactly express the distribution of thermodynamic variations. (These have been derived by many authors now in a wide range of physical settings; see, e.g., Refs. [11] and [12] for lucid reviews.) Second, steady-state thermodynamics (SST) showed that NESSs play a role in nonequilibrium analogous to that of equilibrium macrostates in equilibrium. In this view heat decomposes into the *housekeeping heat* Q_{hk} needed to sustain NESSs and the *excess heat* Q_{ex} dissipated in transitions between them [13–15]. Bolstering SST, recent efforts generalized the Clausius inequality (describing excess heat produced beyond the change in system entropy) to smoothly driven transitions between NESSs [16]. Taken together, these results established an integral fluctuation theorem for the excess work in NESS transitions and, consequently, a generalized Second Law for excess entropy produced beyond housekeeping during driven NESS transitions.

The following extends SST by introducing several new FTs, highlighting correspondences between nonequilibrium and equilibrium relations. First, we provide detailed (i.e., nonintegrated) fluctuation theorems, rather than integral fluctuation theorems for driven NESS transitions. (Integral FTs follow directly, in any case.) This constrains distributions of excess work W_{ex} exerted when controlling nonequilibrium systems. Second, we jointly bound housekeeping and excess work distributions. For example, for time-symmetric driving we show that the joint probability of excess work and housekeeping heat respect the strong constraint:

$$\frac{\Pr(W_{\text{ex}}, Q_{\text{hk}})}{\Pr(-W_{\text{ex}}, -Q_{\text{hk}})} = e^{\beta W_{\text{ex}}} e^{\beta Q_{\text{hk}}}.$$

When the transitions are nonequilibrium excursions between *equilibrium* steady states, this reduces to the Crooks FT [17]. Third, we derive the detailed FTs for entropy production even when temperature varies in space and time. They are expressed in terms of excess environmental entropy production Ω and irreversibility Ψ , even when the irreversibility is “housekeeping” not strictly associated with heat. Finally, we quantify a system’s net *path irreversibility* Ψ with the accumulated violation of detailed balance in the effective dynamic. In the isothermal setting, for example, the irreversibility is the housekeeping heat, maintaining the system in its nonequilib-

rium dynamic: $\Psi = \beta Q_{\text{hk}}$. Importantly, we can determine the minimum housekeeping heat without appealing to the system’s Hamiltonian.

Extending SST in this way reveals universal constraints on excess thermodynamic quantities—effective energies accessible above the housekeeping background. Looking forward, this allows one to analyze nondetailed-balanced stochastic dynamics—and thus contributes an understanding of the role of hierarchy—in the thermodynamics of replication [18] and the thermodynamics of learning [19]. Moreover, this identifies how complex, possibly intelligent, thermodynamic systems leverage (designed or intrinsic) irreversibility in their own state-space to harness energy from structured environments.

B. Synopsis

Section II sets up our approach, introducing notation, discussing input-dependent system dynamics, and establishing fundamental relationships among nonequilibrium thermodynamic quantities. Section III A introduces excess heat and excess work in analogy to classical heat and work. Ultimately though, the related excess environmental entropy production Ω discussed in § III B generalizes these to the case of temperature inhomogeneity over spacetime. Section III C demonstrates that path entropies are the fundamental objects of nonequilibrium thermodynamics. In steady state, unaveraged path entropies reduce to the steady-state surprisal ϕ . Deviations from the asymptotic surprisal contribute to a nonsteady-state additional free energy. All of these quantities play a central role in the subsequent development.

Before delving into irreversibility, though, we first address what is meant by reversibility. Therefore, § IV A and § IV B discuss detailed balance, microscopic reversibility, and the close relationship between them. Section IV C then introduces path dependence and reverse-path dependence and explains how together they yield a system’s irreversibility Ψ .

With this laid out, § V A and § V C derive the detailed FTs in terms of excess environmental entropy production Ω and irreversibility Ψ . One sees that in the isothermal setting $\Psi = \beta Q_{\text{hk}}$ and the excess entropy production is directly related to the excess work. This allows § V E to explain how these results extend SST.

Sections V F and V G finish our investigation of NESS FTs by deriving several integral FTs. This, in effect, extends feedback control, as developed in Refs. [3] and [20], to SST. We note that such environmental feedback is intrinsic to natural systems.

For concreteness, § VI analyzes a simple but biologically important prototype system: voltage-gated sodium

ion channels. These are complex macromolecules that violate detailed balance in order to perform critical biological functioning far from equilibrium. Finally, appendices discuss non-Markovian dynamics and comment on the bounds provided by integral fluctuation theorems for auxiliary variables.

II. DRIVEN STOCHASTIC DYNAMICS

We consider a classical system—the *system under study*—with time-dependent driving via environmentally determined parameters; e.g., time-dependent temperature, voltage, and piston position. Hence, the environmental control input X_t at time t , taking on values $x_t \in \mathcal{X}$, will typically be a vector object. The system under study is assumed to have a countable set \mathcal{S} of states. The random variable \mathcal{S}_t for the state at time t takes on values $s_t \in \mathcal{S}$. We assume that the environment’s control value (current *input*) x and the system’s physical state (current *state*) s are sufficient to determine the system’s net effective energy—the *nonequilibrium potential* $\phi(x, s)$. Even with constant environmental input, the system dynamic need not be detailed balance.

A. Stochastic mesoscopic dynamics and induced state-distributions

We assume the current environmental input x determines the instantaneous stochastic transition dynamic over the system’s observable mesoscopic states. However, that input can itself depend arbitrarily on all previous input and state history. That is, we assume that the \mathcal{S} -to- \mathcal{S} transitions are instantaneously Markovian given the input. Over time, though, different inputs induce different Markov chains over system states.

Note that the Markov assumption is common, although often implicit, and we follow this here to isolate the novel implications of nondetailed-balance dynamics. Nevertheless, the results generalize to infinite Markov order by modeling system states as the observable output of many-to-one mappings of latent states of an input-controllable hidden Markov chain. Appendix A details this generalization.

We do not restrict the environment’s driving process, allowing arbitrary non-Markovity, feedback, and nonstationarity. Thus, the joint system-environment dynamic can be non-Markovian even if the instantaneous system dynamic is. Such a setup is quite general, and so the results to follow extend others known for SST. We also follow stochastic thermodynamics in the use of (arbitrarily small) discrete-time steps. Nevertheless, it is usually

easy to take the continuous-time limit. As, in fact, we do in the example at the end.

Hence, the Markovian dynamic is described by a (possibly infinite) set of input-conditioned transition matrices over the state set \mathcal{S} : $\{\mathsf{T}^{(\mathcal{S} \rightarrow \mathcal{S}|x)}\}_{x \in \mathcal{X}}$, where $\mathsf{T}_{i,j}^{(\mathcal{S} \rightarrow \mathcal{S}|x)} = \Pr(\mathcal{S}_t = s^j | \mathcal{S}_{t-1} = s^i, X_t = x)$ is the probability that the system is in state s^j at time t given that the system was in state s^i at time $t-1$ and the instantaneous environmental input controlling the system was x .

The Perron–Frobenius theorem guarantees that there is a stationary distribution $\boldsymbol{\pi}_x$ over states associated with each fixed input x . These are the state distributions associated with the system’s nonequilibrium steady states (NESSs). For simplicity, and unless otherwise stated, we assume that a fixed input x eventually induces a unique NESS.

We denote distributions over the system states as bold Greek symbols; such as $\boldsymbol{\mu}$. We denote the state random variable \mathcal{S} being distributed according to $\boldsymbol{\mu}$ via $\mathcal{S} \sim \boldsymbol{\mu}$. It will often be convenient to cast $\boldsymbol{\mu}$ as a row-vector, in which case it appears as the bra $\langle \boldsymbol{\mu} |$. Putting this altogether, a sequence of driving inputs updates the state distribution as follows:

$$\begin{aligned} \langle \boldsymbol{\mu}_{t+n} | &= \langle \boldsymbol{\mu}_t | \mathsf{T}^{(\mathcal{S} \rightarrow \mathcal{S}|x_{t:t+n})} \\ &= \langle \boldsymbol{\mu}_t | \mathsf{T}^{(\mathcal{S} \rightarrow \mathcal{S}|x_t)} \mathsf{T}^{(\mathcal{S} \rightarrow \mathcal{S}|x_{t+1})} \dots \mathsf{T}^{(\mathcal{S} \rightarrow \mathcal{S}|x_{t+n-1})} . \end{aligned}$$

(Time indexing here and throughout is denoted by subscripts $t : t'$ that are left-inclusive and right-exclusive.) An infinite driving history $\overleftarrow{x} = \dots x_{-2} x_{-1}$ induces a distribution $\boldsymbol{\mu}(\overleftarrow{x})$ over the state space. The so-called *steady-state distribution* associated with the environmental drive value x , induced by tireless repetition of x , is:

$$\langle \boldsymbol{\pi}_x | = \lim_{n \rightarrow \infty} \langle \boldsymbol{\mu}_0 | \left(\mathsf{T}^{(\mathcal{S} \rightarrow \mathcal{S}|x)} \right)^n .$$

Usefully, $\boldsymbol{\pi}_x$ can also be found as the left eigenvector of $\mathsf{T}^{(\mathcal{S} \rightarrow \mathcal{S}|x)}$ associated with the eigenvalue of unity [21]:

$$\langle \boldsymbol{\pi}_x | = \langle \boldsymbol{\pi}_x | \mathsf{T}^{(\mathcal{S} \rightarrow \mathcal{S}|x)} . \quad (1)$$

The assumption that observable state-to-state transitions are instantaneously Markovian allows the state distribution $\boldsymbol{\mu}$ to summarize the causal relevance of the entire driving history \overleftarrow{x} .

III. ENERGETICS AND ENTROPIES

For later comparison, we recount the basics of a statistical mechanics description of the thermodynamics of a system exchanging energy with a large environment, imposing fixed constraints indexed as x . The

many-body Hamiltonian $\mathcal{H}(x)$ has energy eigenvalues $\{E(x, s)\}$, where s indexes the energy eigenstates. The canonical distribution is $\pi_x(s) = e^{-\beta[E(x, s) - F_{\text{eq}}(x)]}$, at fixed x . This distribution is the equilibrium steady “state” associated with x , where $\beta^{-1} \equiv k_{\text{B}}T$, T is the temperature of the macroscopic environment surrounding the system, and $F_{\text{eq}}(x)$ is the associated equilibrium free energy.

A. Work, heat, and their excesses

Work W is environmentally driven energy change. Within one time-step it is given by [22]:

$$W[x_{n-1} \rightarrow x_n; s_{n-1}] = E(x_n, s_{n-1}) - E(x_{n-1}, s_{n-1}) .$$

Heat Q is the change in system energy due to its internal response to the environmental drive; e.g., a molecule’s change in conformation. Within one time-step the heat is:

$$Q[x_n; s_{n-1} \rightarrow s_n] = E(x_n, s_n) - E(x_n, s_{n-1}) .$$

Over the course of driving the system from $t = 0$ to $t = N\Delta t = \tau$, the net energy change is then:

$$\begin{aligned} \Delta E &= E(x_N, s_N) - E(x_0, s_0) \\ &= W + Q , \end{aligned}$$

where the net work and net heat are:

$$W = \sum_{n=1}^N W[x_{n-1} \rightarrow x_n; s_{n-1}]$$

and:

$$Q = \sum_{n=1}^N Q[x_n; s_{n-1} \rightarrow s_n] ,$$

respectively. Here, and later on, Δ applied to a quantity refers to its change over one time step, where the step is given by context.

When the system strongly couples to a substrate with uncontrolled energy fluxes, steady-state dynamics are often established far from equilibrium, even when environmental parameters are held fixed. That is, for fixed driving $\dots xxxxx \dots$, the system settles down to a NESS with a distribution over observable system states given by the *nonequilibrium potential* $\phi(x, s)$:

$$\pi_x(s) = e^{-\phi(x, s)} . \quad (2)$$

In this, $\phi(x, s)$ plays a role roughly analogous to en-

ergy eigenvalues. Thus, the thermodynamics of accessible energetics—the excess heat generated and work irretrievably performed in driving between NESSs—follows analogously to its equilibrium counterpart. This is complementary to recent SST studies [14–16, 23].

If steady-state free energies $F_{\text{ss}}(x)$ and *effective* energies $E_{\text{eff}}(x, s)$ could be uniquely (and usefully) defined, then the nonequilibrium potential would be:

$$\phi(x, s) = \beta[E_{\text{eff}}(x, s) - F_{\text{ss}}(x)] .$$

However, the assignment of steady-state free energies is problematic. Nevertheless, $\phi(x, s)$ retains meaning since it quantifies the steady-state *surprisal* of observing state s :

$$\phi(x, s) = -\ln \pi_x(s) .$$

The surprisal is Shannon’s *self-information* [24]—the un-averaged individual-event entropy measuring how surprising a specific event is. Intuitively, we must do work to make otherwise unlikely things happen.

SST’s *excess* work and heat can be defined via changes in steady-state surprisal ϕ , analogous to how equilibrium quantities are in terms of energy changes. For clarity, we temporarily restrict ourselves to the isothermal setting, but we can easily adapt to time-varying temperatures.

Excess work W_{ex} is environmentally driven change in nonequilibrium potential:

$$\begin{aligned} W_{\text{ex}}[x_{n-1} \rightarrow x_n; s_{n-1}] \\ = \beta^{-1}[\phi(x_n, s_{n-1}) - \phi(x_{n-1}, s_{n-1})] , \end{aligned}$$

over one time-step. *Excess heat* Q_{ex} is the change in nonequilibrium potential due to the system’s response:

$$Q_{\text{ex}}[x_n; s_{n-1} \rightarrow s_n] = \beta^{-1}[\phi(x_n, s_n) - \phi(x_n, s_{n-1})] ,$$

over one time-step. When driving from $t = 0$ to $t = N\Delta t = \tau$, the net change in nonequilibrium potential is:

$$\begin{aligned} \Delta\phi &= \phi(x_N, s_N) - \phi(x_0, s_0) \\ &= \beta(W_{\text{ex}} + Q_{\text{ex}}) \\ &= -\ln \frac{\pi_{x_N}(s_N)}{\pi_{x_0}(s_0)} , \end{aligned} \quad (3)$$

where the net excess work and net excess heat are:

$$W_{\text{ex}} = \sum_{n=1}^N W_{\text{ex}}[x_{n-1} \rightarrow x_n; s_{n-1}] \quad (4)$$

and:

$$Q_{\text{ex}} = \sum_{n=1}^N Q_{\text{ex}}[x_n; s_{n-1} \rightarrow s_n], \quad (5)$$

respectively.

This approach to excess heat Q_{ex} coincides with SST's definition and reduces to total heat in equilibrium transitions. Importantly, it follows as closely as possible the equilibrium approach to total heat Q outlined above and deviates from the typical starting point: $Q_{\text{ex}} \equiv Q - Q_{\text{hk}}$, where Q_{hk} is the so-called *housekeeping heat*. In contrast, excess work W_{ex} does *not* reduce to the total work in equilibrium transitions. Rather, W_{ex} goes over to $W - \Delta F_{\text{eq}}$, if the steady states are near equilibrium. And, this fortuitously coincides with its previous narrower use in describing transitions atop equilibrium steady states—the work exerted beyond the change in free energy [25].

The excess heat Q_{ex} can be interpreted as the heat dissipated during transitions between NESSs. Similarly, the excess work W_{ex} can be interpreted as the work that *would* be dissipated if the system is allowed to relax back to a NESS. The difference between excess work W_{ex} and *dissipated work*, denoted W_{diss} , depends on a notion of excess nonequilibrium free energy, discussed shortly.

This framing leads us to see that heat is how small, possibly intelligent, systems store and transform energy via their *own agency*. This stance also moves us away from any unjustified biases that heat is necessarily wasteful. For example, an increase in heat may indicate that a system has harvested energy, and the emission of heat may indicate an intrinsic computation [26] in the system's state-space. The *efficiency* of the tradeoff—spending stored energy to achieve some utility—then comes into question. It is inefficiency in this sense that is by its nature wasteful.

B. Excess environmental entropy production

In isothermal transitions between equilibrium steady states, the environmental entropy production is [17]:

$$\begin{aligned} \Omega_{\text{eq}} &= \beta(W - \Delta F_{\text{eq}}) \\ &= -\beta Q - \ln \frac{\pi_{x_N}(s_N)}{\pi_{x_0}(s_0)}. \end{aligned}$$

This extends to SST by defining the *excess environmental entropy production*:

$$\begin{aligned} \Omega &= \beta W_{\text{ex}} \\ &= -\beta Q_{\text{ex}} - \ln \frac{\pi_{x_N}(s_N)}{\pi_{x_0}(s_0)}, \end{aligned} \quad (6)$$

This has also been referred to as the “nonadiabatic component of entropy production” [16, 23, 27, 28]. Note that $-\ln(\pi_{x_N}(s_N)/\pi_{x_0}(s_0)) = \Delta\phi$, recovering Eq. (3)'s change in nonequilibrium potential ϕ .

Recalling the definitions of Q_{ex} in terms of steady state surprisals and $\phi(x, s) = -\ln \pi_x(s)$, we see that:

$$\begin{aligned} e^{-\beta Q_{\text{ex}}} &= e^{-\sum_{n=1}^N [\phi(x_n, s_n) - \phi(x_n, s_{n-1})]} \\ &= \prod_{n=1}^N \frac{\pi_{x_n}(s_n)}{\pi_{x_n}(s_{n-1})}. \end{aligned} \quad (7)$$

And so, Eq. (6) gives:

$$\begin{aligned} e^{\Omega(x_0:N+1, s_0:N)} &= \frac{\pi_{x_0}(s_0)}{\pi_{x_N}(s_N)} \prod_{n=1}^N \frac{\pi_{x_n}(s_n)}{\pi_{x_n}(s_{n-1})} \\ &= \prod_{n=0}^{N-1} \frac{\pi_{x_n}(s_n)}{\pi_{x_{n+1}}(s_n)}. \end{aligned} \quad (8)$$

If temperature varies, then the above still holds if we replace the equilibrium probabilities with the temperature-dependent equilibrium probabilities. Thus, to go beyond the isothermal setting, we use Eq. (8) as the defining relationship for the excess environmental entropy production Ω . If temperature is spatially homogeneous, then it is equivalent to:

$$\Omega = \Delta\phi - \frac{1}{k_B} \int \frac{\delta Q_{\text{ex}}}{T}.$$

However, spatially inhomogeneous temperatures can also be addressed by folding temperature dependence into the environmental input x .

We return to these expressions and explore their role in generalized fluctuation theorems once we develop the necessary quantitative notions of irreversibility.

C. Path entropies

In steady state, the system state probability distribution has a Boltzmann exponential dependence on the effective energies. Naturally, out of steady state the distribution is something different. There is a nonsteady-state free energy associated with this out-of-steady-state distribution, since the system can do work (or computations) at the cost of relaxing the distribution.

Nonsteady-state free energies are controlled by path entropies, which come in several varieties. Here, we are especially interested in the controllable unaveraged state surprisals induced by the driving path \overleftarrow{x} :

$$h^{(s|x_{-\infty:t+1})} = -\ln \Pr(\mathcal{S}_t = s | x_{-\infty:t+1}). \quad (9)$$

Since a semi-infinite history induces a particular distribution over system states, this can be usefully recast in terms of the initial distribution $\boldsymbol{\mu}_0$ induced by the path $x_{-\infty:1}$ and the driving history $x_{1:t+1}$ since then:

$$\begin{aligned} h^{(s|\boldsymbol{\mu}_0, x_{1:t+1})} &= -\ln \Pr(\mathcal{S}_t = s | \mathcal{S}_0 \sim \boldsymbol{\mu}_0, x_{1:t+1}) \quad (10) \\ &= -\ln \langle \boldsymbol{\mu}_0 | \mathsf{T}^{(\mathcal{S} \rightarrow \mathcal{S} | x_{1:t+1})} | s \rangle, \end{aligned}$$

where $\Pr(\mathcal{S}_t = s | \mathcal{S}_0 \sim \boldsymbol{\mu}_0, x_{1:t+1})$ is the probability that the state is s at time t , under the measure induced when the initial state $\mathcal{S}_0 \sim \boldsymbol{\mu}_0$ (distributed according to $\boldsymbol{\mu}_0$) [29] and given the driving history $x_{1:t+1} = x_1 \dots x_t$ since the initial time.

Alternatively, consider the distribution $\boldsymbol{\mu}$ induced from a start distribution by the driving history since the start. Then the path-induced state-surprisal can be expressed simply in terms of the *present* environmental-history-induced distribution over system states and the candidate state s :

$$\begin{aligned} h^{(s|\boldsymbol{\mu})} &= -\ln \Pr(\mathcal{S}_t = s | \mathcal{S}_t \sim \boldsymbol{\mu}) \quad (11) \\ &= -\ln \langle \boldsymbol{\mu} | s \rangle. \end{aligned}$$

Thermodynamic units of entropy are recovered by multiplying the Shannon-like path surprisals by Boltzmann's constant: $\boldsymbol{s} = k_B h$.

Averaging the path-induced state-surprisal over states gives a genuine input-conditioned Shannon entropy:

$$\begin{aligned} \langle h^{(s_t | \overleftarrow{x}_t)} \rangle_{\Pr(s_t | \overleftarrow{x}_t)} &= -\left\langle \ln \Pr(s_t | \overleftarrow{x}_t) \right\rangle_{\Pr(s_t | \overleftarrow{x}_t)} \\ &= -\sum_{s_t} \Pr(s_t | \overleftarrow{x}_t) \ln \Pr(s_t | \overleftarrow{x}_t) \\ &= \mathsf{H}[\mathcal{S}_t | \overleftarrow{X}_t = \overleftarrow{x}_t], \quad (12) \end{aligned}$$

where $\mathsf{H}[\cdot | \cdot]$ is the conditional Shannon entropy in units of nats.

It follows directly that the state-averaged path entropy $k_B \mathsf{H}[\mathcal{S}_t | \overleftarrow{x}_t]$ is an extension of the system's steady-state nonequilibrium entropy S_{ss} . In steady-state, the state-averaged path entropy reduces to:

$$\begin{aligned} k_B \mathsf{H}[\mathcal{S}_t | \overleftarrow{X}_t = \dots xxx] &= -k_B \mathsf{H}[\mathcal{S}_t | \mathcal{S}_t \sim \boldsymbol{\pi}_x] \\ &= -k_B \sum_{s \in \mathcal{S}} \pi_x(s) \ln \pi_x(s) \quad (13) \\ &= S_{\text{ss}}(x). \end{aligned}$$

The system steady-state nonequilibrium entropy S_{ss} has been discussed as a fundamental entity in SST; e.g., see Refs. [14] and [16]. However, Eq. (12) ($\times k_B$) gives the appropriate extension for the thermodynamic entropy of a nonequilibrium system that is *not in steady state*. Rather, it is the entropy over system states given the

entire history of environmental driving.

When \mathcal{S} is the set of microstates, rather than, say, observable mesoscopic states, the unaveraged nonequilibrium free energy F enjoys the familiar relationship between energy E and (path) entropy \boldsymbol{s} :

$$F^{(s_t | x_{-\infty:t+1})} \equiv E(x_t, s_t) - T \boldsymbol{s}^{(s_t | x_{-\infty:t+1})} \quad (14)$$

$$= F_{\text{eq}}(x_t) + \beta^{-1} \ln \frac{\Pr(s_t | x_{-\infty:t+1})}{\pi_{x_t}(s_t)}. \quad (15)$$

Or, averaging over states:

$$\begin{aligned} \mathcal{F}(t) &= U(t) - \beta^{-1} \mathsf{H}[\mathcal{S}_t | x_{-\infty:t+1}] \quad (16) \\ &= F_{\text{eq}}(x_t) + \beta^{-1} D_{\text{KL}}(\Pr(\mathcal{S}_t | x_{-\infty:t+1}) || \boldsymbol{\pi}_{x_t}), \end{aligned}$$

where $\mathcal{F}(t)$ is the expected instantaneous nonequilibrium free energy, $U(t)$ is the expected instantaneous thermal energy, and $D_{\text{KL}}(\cdot || \cdot)$ is the Kullback–Leibler divergence [24]. Recognizing $k_B \mathsf{H}[\mathcal{S}_t | x_{-\infty:t+1}]$ as the natural extension of a system's thermodynamic entropy, Eq. (16) is familiar from equilibrium thermodynamics, but it is now applicable arbitrarily far from equilibrium and at any time t using the instantaneous temperature. This is not the first statement of such a generalized relationship; compare, e.g., Refs. [30, 31]. In equilibrium, the expected value of the path entropy (using microstates) reduces to *the* equilibrium entropy of a system.

In the setting of effective states and NESS surprisals, we can no longer directly use Eq. (14). Nevertheless, by analogy with Eq. (15), we can still identify the *nonsteady-state addition* $\gamma(\cdot | \cdot)$ to free energy as:

$$\beta^{-1} \gamma(s | \boldsymbol{\mu}, x) \equiv \beta^{-1} \ln \frac{\Pr(\mathcal{S}_t = s | \mathcal{S}_{t-1} \sim \boldsymbol{\mu}, X_t = x)}{\pi_x(s)}.$$

Expressed differently, it is:

$$\begin{aligned} \gamma(s | \boldsymbol{\mu}, x) &= h^{(s|\boldsymbol{\pi}_x)} - h^{(s|\boldsymbol{\mu}, x)} \\ &= \phi(x, s) - h^{(s|\boldsymbol{\mu}, x)}. \end{aligned}$$

Averaging over states this becomes the Kullback–Leibler divergence between nonsteady state and steady state distributions:

$$\langle \gamma(s | \boldsymbol{\mu}, x) \rangle = D_{\text{KL}}[\Pr(\mathcal{S}_t | \mathcal{S}_{t-1} \sim \boldsymbol{\mu}, X_t = x) || \boldsymbol{\pi}_x],$$

which is nonnegative.

Identifying the nonsteady-state contribution to the free energy allows us to introduce the *dissipated work*:

$$W_{\text{diss}} \equiv W_{\text{ex}} - \beta^{-1} \Delta \gamma,$$

to account for the fact that excess work is not fully dissipated until the distribution relaxes back to steady state

π_x . An important consequence is that the excess work dissipated can be reclaimed by a subsequent “fluctuation” with $W_{\text{ex}} < 0$ in the midst of a driven nonequilibrium excursion.

The role of the nonsteady-state contribution to free energy will be apparent in the FTs to come shortly. This generalizes similar FTs that are restricted to starting and possibly ending in a steady state π_x . The generalization here is key to analyzing complex systems, since many simply cannot be initiated in a steady state without losing their essential character.

IV. IRREVERSIBILITY

To emphasize, the preceding did not reference and does not require detailed balance. However, to ground the coming development, we need to describe the roles of reversibility, detailed balance, and their violations. At a minimum, this is due to most FTs assuming reversibility of the effective dynamic over states. Having established the necessary concepts and giving a measure of the irreversibility of the effective dynamic, we finally move on to FTs for nondetailed balanced processes.

A. Detailed balance

Transitioning from state a to state b , say, invoking detailed balance assumes that:

$$\frac{\Pr(\mathcal{S}_n = a | \mathcal{S}_{n-1} = b, X_n = x)}{\Pr(\mathcal{S}_n = b | \mathcal{S}_{n-1} = a, X_n = x)} = \frac{\pi_x(a)}{\pi_x(b)}.$$

Though, we do *not* assume detailed balance over the states considered here, we refer to it occasionally. For example, assuming detailed balance, microscopic reversibility and the standard Crooks fluctuation theorem (CFT) follow almost immediately.

In contrast, complex systems sustained out of equilibrium by an active substrate generically evolve via nondetailed-balance dynamics. To wit, many examples of nondetailed-balance dynamics are exhibited by chemical kinetics in biological systems [32–34]. We conclude with a thorough-going thermodynamic analysis of one neurobiological example.

B. Microscopic reversibility

Consider a particular realization of interleaved environmental-input sequence $\dots x^1 x^2 \dots x^{N-1} \dots$ and

system-state sequence $\dots s^1 s^2 \dots s^{N-1} \dots$:

$$\begin{array}{ccccccc} & & & \mathbf{x}, \mathbf{x}^{\text{R}} & & & \\ & & & \overbrace{\hspace{10em}} & & & \\ x^0 & x^1 & x^2 & \dots & x^{N-1} & x^N & \\ & \underbrace{\hspace{10em}}_{\mathbf{s}} & & & & & \\ s^0 & s^1 & s^2 & \dots & s^{N-2} & s^{N-1} & s^N \\ & & & \underbrace{\hspace{10em}}_{\mathbf{s}_{\leftarrow}^{\text{R}}} & & & \end{array}$$

There are several length- $(N - 1)$ subsequences in play here, including the forward trajectory $\mathbf{x} = x^1 x^2 \dots x^{N-2} x^{N-1}$ of the environmental driving and the forward trajectory $\mathbf{s} = s^1 s^2 \dots s^{N-2} s^{N-1}$ of the state sequence. Furthermore, let $\mathbf{x}^{\text{R}} = x^{N-1} x^{N-2} \dots x^2 x^1$ be the time-reversal of the environmental driving \mathbf{x} and $\mathbf{s}_{\leftarrow}^{\text{R}} = s^{N-2} s^{N-3} \dots s^1 s^0$ the time-reversal of the time-shifted state sequence \mathbf{s} .

For example, if $\mathcal{X} = \{0, 1\}$ and $\mathcal{S} = \{a, b, c\}$, then \mathbf{x} may be the sequence 00101110...11000010 and \mathbf{s} the sequence *acaaaaba...abaccabc*. Then \mathbf{x}^{R} is the sequence 01000011...01110100. Taking the time reversal of the state sequence, we have *cbaccaba...abaaaaca*. However, since $\mathbf{s}_{\leftarrow}^{\text{R}}$ is also time-shifted by one time-step, we must drop the first c and append another symbol, say a . Then $\mathbf{s}_{\leftarrow}^{\text{R}}$ is the sequence *baccaba...abaaaaca*.

Let Q_{F} be the excess heat of the joint forward sequences \mathbf{x} and $s^0 \mathbf{s}$, according to Eq. (5). By definition, a system–environment effective dynamic is *microscopically reversible* if:

$$\frac{\Pr(\mathcal{S}_{1:N} = \mathbf{s} | \mathcal{S}_0 = s^0, X_{1:N} = \mathbf{x})}{\Pr(\mathcal{S}_{1:N} = \mathbf{s}_{\leftarrow}^{\text{R}} | \mathcal{S}_0 = s^{N-1}, X_{1:N} = \mathbf{x}^{\text{R}})} = e^{-\beta Q_{\text{F}}},$$

for any $s^0 \in \mathcal{S}$, $\mathbf{s} \in \mathcal{S}^{N-1}$, and $\mathbf{x} \in \mathcal{X}^{N-1}$. As a useful visual aid, we can re-express this as:

$$\frac{\Pr(s^0 \xrightarrow{x^1} s^1 \dots s^{N-2} \xrightarrow{x^{N-1}} s^{N-1} | s^0, \mathbf{x})}{\Pr(s^0 \xleftarrow{x^1} s^1 \dots s^{N-2} \xleftarrow{x^{N-1}} s^{N-1} | s^{N-1}, \mathbf{x}^{\text{R}})} = e^{-\beta Q_{\text{F}}}.$$

Otherwise, microscopic reversibility is broken.

Although, microscopic reversibility has also been referred to as a “detailed fluctuation theorem”, it is actually an assumption appropriate only in special cases. For example, Eq. (7) shows that if the dynamics are Markovian over states (given input) and obey detailed balance (à la §IV A), then microscopic reversibility is satisfied for arbitrary non-Markovian inputs. In essence, this is the justification of microscopic reversibility suggested by Crooks [2, 17] from which his eponymous fluctuation theorem follows.

In this view, detailed balance and microscopic reversibility are effectively the same assumption since each implies the other. Section § V C generalizes the CFT to describe fluctuation laws in the absence of microscopic

reversibility.

C. Path dependence and irreversibility

The importance of state-space path dependence is captured via an informational quantity Υ we call the *path relevance* of a state sequence $s_{1:N}$ given initial state s_0 and input sequence $x_{1:N}$:

$$\Upsilon(s_{1:N}|s_0, x_{1:N}) \equiv \ln \frac{\Pr(s_{1:N}|s_0, x_{1:N})}{\prod_{n=1}^{N-1} \pi_{x_n}(s_n)}. \quad (17)$$

(The branching Pythagorean letter Υ recognizes its ancient symbolism—divergent consequences of choosing one path over another.) Note that the equilibrium probabilities in the denominator do not depend on the original state, whereas the numerator (even after factoring) depends on state-to-state transitions. As the environmental input drives the system probability density through its state-space, path relevance develops in the state sequence. A joint sequence lacks path relevance, if $\Upsilon = 0$ for that sequence.

Whenever state transitions are Markovian given the input, the numerator in Eq. (17) simplifies to:

$$\Pr(s_{1:N}|x_{1:N}, s_0) = \prod_{n=1}^{N-1} \Pr(s_n|s_{n-1}, x_n),$$

and the path relevance becomes:

$$\Upsilon = \sum_{n=1}^{N-1} \ln \frac{\Pr(s_n|s_{n-1}, x_n)}{\pi_{x_n}(s_n)}.$$

Thus, there is path relevance even for Markov processes. The actual driving history matters. When a system is *non*-Markovian, there are yet additional contributions to path relevance.

Path relevance of a particular state sequence given a particular driving is a system feature, regardless of the environment in which the system finds itself. However, expectation values involving the above relationship can reflect the environment's nature.

For our development, we find it useful to consider both the forward-path dependence and the reverse-path dependence of a particular joint sequence: $x^1 \dots x^{N-1}$ and

$s^0 \dots s^{N-1}$. The *forward-path dependence* is as expected:

$$\begin{aligned} \Upsilon &= \Upsilon(\mathbf{s}|s^0, \mathbf{x}) \\ &= \ln \frac{\Pr(\mathbf{s}|s^0, \mathbf{x})}{\prod_{n=1}^{N-1} \pi_{x_n}(s^n)}, \end{aligned} \quad (18)$$

and, similarly, the *reverse-path dependence* is:

$$\begin{aligned} \mathcal{J} &= \mathcal{J}(\mathbf{s}|s^0, \mathbf{x}) \\ &= \Upsilon(\mathbf{s}_{\leftarrow}^R|s^{N-1}, \mathbf{x}^R) \\ &= \ln \frac{\Pr(\mathbf{s}_{\leftarrow}^R|s^{N-1}, \mathbf{x}^R)}{\prod_{n=1}^{N-1} \pi_{x^n}(s^{n-1})}. \end{aligned} \quad (19)$$

And, finally, we have the *irreversibility*:

$$\Psi \equiv \Upsilon - \mathcal{J}, \quad (20)$$

the *net directional relevance*—of a particular path \mathbf{s} given s^0 and \mathbf{x} . Nonzero Ψ indicates the irrevocable consequences of path traversal. *Microscopically reversible dynamics have $\Psi = 0$ for all paths with nonzero probability*, indicating no divergence in path branching anywhere through the state-space. And so, $\Psi = 0$ for all paths with nonzero probability if and only if the dynamic satisfies detailed balance. In short, Ψ quantifies the imbalance in path reciprocity along a driven state-sequence.

Sometimes \mathcal{J} can be $-\infty$ for an allowed forward path $\Pr(\mathbf{s}|s^0, \mathbf{x}) > 0$, corresponding to a forbidden reverse path $\Pr(\mathbf{s}_{\leftarrow}^R|s^{N-1}, \mathbf{x}^R) = 0$. This is a situation that never arises with detailed balance dynamics. Such paths are infinitely irreversible: $\Psi = \infty$.

V. GENERALIZED FLUCTUATION THEOREMS FOR NONEQUILIBRIUM SYSTEMS

Absent microscopic reversibility, the architecture of transitions over state-space matters. More concretely, we will constructively show how this architecture affects the nonequilibrium thermodynamics of complex systems.

A. Generalized Detailed Fluctuation Theorem

Assume the system under study starts from some distribution μ_F and that the associated reverse trajectory (when starting from some other distribution μ_R) is allowed—that is, it has nonzero probability. Then the ratio of conditional probabilities of a state sequence (given a driving sequence) to the reversed state sequence (given reversed driving) is:

$$\begin{aligned}
& \frac{\Pr(\boldsymbol{\mu}_F \xrightarrow{x^0} s^0 \xrightarrow{x^1} s^1 \dots s^{N-2} \xrightarrow{x^{N-1}} s^{N-1} | \boldsymbol{\mu}_F, x^0 \mathbf{x})}{\Pr(s^0 \xleftarrow{x^1} s^1 \dots s^{N-2} \xleftarrow{x^{N-1}} s^{N-1} | \boldsymbol{\mu}_R | \boldsymbol{\mu}_R, x^N \mathbf{x}^R)} \\
&= \frac{\Pr(\mathcal{S}_{0:N} = s^0 \mathbf{s} | \mathcal{S}_{-1} \sim \boldsymbol{\mu}_F, X_{0:N} = x^0 \mathbf{x})}{\Pr(\mathcal{S}_{0:N} = s^{N-1} \mathbf{s}_{\leftarrow}^R | \mathcal{S}_{-1} \sim \boldsymbol{\mu}_R, X_{0:N} = x^N \mathbf{x}^R)} \\
&= \frac{\Pr(s^0 | \boldsymbol{\mu}_F, x^0)}{\Pr(s^{N-1} | \boldsymbol{\mu}_R, x^N)} \frac{\Pr(\mathbf{s} | s^0, \mathbf{x})}{\Pr(\mathbf{s}_{\leftarrow}^R | s^{N-1}, \mathbf{x}^R)} \\
&= \frac{\Pr(s^0 | \boldsymbol{\mu}_F, x^0)}{\pi_{x^0}(s^0)} \frac{\pi_{x^N}(s^{N-1})}{\Pr(s^{N-1} | \boldsymbol{\mu}_R, x^N)} \frac{\Pr(\mathbf{s} | s^0, \mathbf{x})}{\prod_{n=1}^{N-1} \pi_{x^n}(s^n)} \frac{\prod_{n=1}^{N-1} \pi_{x^n}(s^{n-1})}{\Pr(\mathbf{s}_{\leftarrow}^R | s^{N-1}, \mathbf{x}^R)} \prod_{n=0}^{N-1} \frac{\pi_{x^n}(s^n)}{\pi_{x^{n+1}}(s^n)} \\
&= \frac{\Pr(s^0 | \boldsymbol{\mu}_F, x^0)}{\pi_{x^0}(s^0)} \frac{\pi_{x^N}(s^{N-1})}{\Pr(s^{N-1} | \boldsymbol{\mu}_R, x^N)} e^{\Psi_F} e^{\Omega_F} \\
&= e^{\gamma(s^0 | \boldsymbol{\mu}_F, x^0) - \gamma(s^{N-1} | \boldsymbol{\mu}_R, x^N)} e^{\Omega_F + \Psi_F}, \tag{21}
\end{aligned}$$

where $\Omega_F = \Omega(X_{0:N+1} = x^0 \mathbf{x}^N, \mathcal{S}_{0:N} = s^0 \mathbf{s})$ is the excess environmental entropy production in the forward trajectory, $\Psi_F = \Psi(\mathcal{S}_{0:N} = s^0 \mathbf{s} | X_{1:N} = \mathbf{x})$ is the irreversibility of the forward trajectory, and $\gamma(s | \boldsymbol{\mu}, x) = \ln(\Pr(s | \boldsymbol{\mu}, x) / \pi_x(s))$ is the nonsteady-state addition to free energy associated with being in the nonsteady-state distribution $\boldsymbol{\mu}$ with environmental drive x .

Since Ψ_F can diverge for forward paths with nonzero probability, we typically rewrite Eq. (21) as the “less divergent” expression:

$$\begin{aligned}
& e^{-\gamma_F} \Pr(s^0 \mathbf{s} | \boldsymbol{\mu}_F, x^0 \mathbf{x}) e^{-\Psi_F} \\
&= e^{-\gamma_R} \Pr(s^{N-1} \mathbf{s}_{\leftarrow}^R | \boldsymbol{\mu}_R, x^N \mathbf{x}^R) e^{\Omega_F}. \tag{22}
\end{aligned}$$

Eq. (22) is the fundamental relation for all that follows: it relates the probabilities of forward and reverse trajectories via entropy production Ω_F of the forward path, irreversibility Ψ_F of the forward path, and change $\beta^{-1}(\gamma_F - \gamma_R)$ in the nonsteady-state addition to free energy between the forward and reverse start-distributions.

In what follows it will be all too easy to write seemingly divergent expressions. Such divergences do not manifest themselves when taking expectation values for physical quantities involving them, since they come weighted with zero probability. This is similar to the reasonable convention for Shannon entropies that $0 \log 0 = 0$. Nevertheless, caution is advised when probabilities vanish.

B. Simplifications

Before proceeding and to aid understanding, let’s consider several special cases. If the forward drive or protocol begins with the system equilibrated to the static environmental drive x^0 , then $\boldsymbol{\mu}_F = \boldsymbol{\pi}_{x^0}$ and $\gamma(s^0 | \boldsymbol{\mu}_F, x^0) = 0$.

Similarly, if the reverse protocol begins with the system equilibrated to the static environmental drive x^N , then $\boldsymbol{\mu}_R = \boldsymbol{\pi}_{x^N}$ and $\gamma(s^{N-1} | \boldsymbol{\mu}_R, x^N) = 0$. In this case, Eq. (22) simplifies to:

$$\frac{\Pr(\boldsymbol{\pi}_{x^0} \xrightarrow{x^0} s^0 \dots \xrightarrow{x^{N-1}} s^{N-1} | \boldsymbol{\pi}_{x^0}, x^0 \mathbf{x})}{\Pr(s^0 \xleftarrow{x^1} \dots s^{N-1} \xleftarrow{x^N} \boldsymbol{\pi}_{x^N} | \boldsymbol{\pi}_{x^N}, x^N \mathbf{x}^R)} = e^{\Omega_F + \Psi_F}.$$

As a separate matter, if the dynamics are microscopically reversible, then $\Psi = 0$. Consider the very special case where (i) the dynamics are microscopically reversible, (ii) the forward driving begins with the system equilibrated with x^0 , and (iii) the reverse driving begins with the system equilibrated with x^N . Then, the ratio of probabilities of observing a forward state sequence (given forward driving) and observing the reversal of that state sequence (given the reversal of that driving) is simply e^{Ω_F} . That is, the forward sequence is exponentially more likely if it has positive entropy production.

Apparently, the more general case is more nuanced and, beyond depending on a nonsteady-state starting distribution, it depends strongly on the architecture of branching transitions among states.

Another interesting special case is if $x^N = x^{N-1}$ and $\boldsymbol{\mu}_R$ is the distribution that the forward driving induces from $\boldsymbol{\mu}_F$. Then the *dissipated work* $W_{\text{diss}} \equiv W_{\text{ex}} - \beta^{-1} \Delta \gamma$ associated with the forward trajectory comes into play. (Recall that $\beta^{-1} \Delta \gamma$ is the change in nonsteady-state contributions to free energy.) Then the ratio of forward- and

reverse-path probabilities is:

$$\begin{aligned} & \frac{\Pr(\boldsymbol{\mu}_F \xrightarrow{x^0} \dots \xrightarrow{x^{N-1}} s^{N-1} | \boldsymbol{\mu}_F, \mathbf{x})}{\Pr(s^0 \xleftarrow{x^1} \dots \xleftarrow{x^N} \boldsymbol{\mu}(\boldsymbol{\mu}_F, \mathbf{x}) | \boldsymbol{\mu}(\boldsymbol{\mu}_F, \mathbf{x}), \mathbf{x}^R)} \\ &= e^{\Psi_F} e^{\beta[W_{\text{ex}} - \beta^{-1} \Delta\gamma]} \\ &= e^{\Psi_F} e^{\beta W_{\text{diss}}} . \end{aligned}$$

Even in the case of microscopic reversibility, this is useful since it generalizes previous FTs to nonequilibrium start and end distributions. In the case of microscopic reversibility, $\Psi_F = 0$ and so the ratio of forward- and reverse-path probabilities from any nonequilibrium start and end distribution is exponential $e^{\beta W_{\text{diss}}}$ in the dissipated work. Thus, an experimental test of this result is one with time-symmetric driving. The forward protocol corresponds to the first half of the driving while the reverse protocol is the second half. Clearly, the final nonequilibrium distribution for the forward protocol is the same as the initial nonequilibrium distribution for the reverse protocol. The dissipated work then corresponds to that dissipated in the first half of the driving. Practically, in cases where the dynamic is not microscopically reversible, this allows experimentally extracting the system's irreversibility Ψ .

C. Generalized Crooks Fluctuation Theorem

We can now turn to the irreversible analog of the Crooks Fluctuation Theorem (CFT).

First, we note that both entropy production Ω and irreversibility Ψ are odd under time reversal. Explicitly, we have:

$$\begin{aligned} & \Omega(X_{0:N+1} = x^0 \mathbf{x}^N, S_{0:N} = s^0 \mathbf{s}) \\ &= \ln \prod_{n=0}^{N-1} \frac{\pi_{x^{n+1}}(s^n)}{\pi_{x^n}(s^n)} \\ &= -\ln \prod_{n=0}^{N-1} \frac{\pi_{x^{n+1}}(s^n)}{\pi_{x^n}(s^n)} \\ &= -\ln \prod_{n=0}^{N-1} \frac{\pi_{x^{N-n}}(s^{N-1-n})}{\pi_{x^{N-1-n}}(s^{N-1-n})} \\ &= -\Omega(X_{0:N+1} = x^N \mathbf{x}^R x^0, S_{0:N} = s^{N-1} \mathbf{s}_{\leftarrow}^R) \end{aligned}$$

and:

$$\begin{aligned} & \Psi(S_{0:N} = s^0 \mathbf{s} | X_{1:N} = \mathbf{x}) \\ &= \ln \left[\frac{\Pr(\mathbf{s} | s^0, \mathbf{x})}{\Pr(\mathbf{s}_{\leftarrow}^R | s^{N-1}, \mathbf{x}^R)} \prod_{n=1}^{N-1} \frac{\pi_{x^n}(s^{n-1})}{\pi_{x^n}(s^n)} \right] \\ &= -\ln \left[\frac{\Pr(\mathbf{s}_{\leftarrow}^R | s^{N-1}, \mathbf{x}^R)}{\Pr(\mathbf{s} | s^0, \mathbf{x})} \prod_{n=1}^{N-1} \frac{\pi_{x^n}(s^n)}{\pi_{x^n}(s^{n-1})} \right] \\ &= -\Psi(S_{0:N} = s^{N-1} \mathbf{s}_{\leftarrow}^R | X_{1:N} = \mathbf{x}^R) . \end{aligned}$$

For brevity, let $\Omega_F \equiv \Omega(X_{0:N+1} = x^0 \mathbf{x}^N, S_{0:N} = s^0 \mathbf{s})$ and $\Omega_R \equiv \Omega(X_{0:N+1} = x^N \mathbf{x}^R x^0, S_{0:N} = s^{N-1} \mathbf{s}_{\leftarrow}^R)$. And, similarly, $\Psi_F \equiv \Psi(S_{0:N} = s^0 \mathbf{s} | X_{1:N} = \mathbf{x})$ and $\Psi_R \equiv \Psi(S_{0:N} = s^{N-1} \mathbf{s}_{\leftarrow}^R | X_{1:N} = \mathbf{x}^R)$. In this notation, we just established that $\Omega_F = -\Omega_R$ and $\Psi_F = -\Psi_R$.

Second, if we now choose $\boldsymbol{\mu}_F = \boldsymbol{\pi}_{\mathbf{x}^0}$ and $\boldsymbol{\mu}_R = \boldsymbol{\pi}_{\mathbf{x}^N}$ and marginalize over all possible state trajectories, we find that the joint probability of entropy production and irreversibility given the driving protocol starting from an equilibrium distribution is:

$$\begin{aligned} & \Pr(\Omega, \Psi | \boldsymbol{\pi}_{\mathbf{x}^0}, x^0 \mathbf{x}^N) \\ &= \sum_{s^0 \mathbf{s} \in \mathcal{S}^N} \Pr(s^0 \mathbf{s} | \boldsymbol{\pi}_{\mathbf{x}^0}, x^0 \mathbf{x}) \delta_{\Omega, \Omega_F} \delta_{\Psi, \Psi_F} \\ &= \sum_{s^0 \mathbf{s} \in \mathcal{S}^N} e^{\Omega_F} e^{\Psi_F} \Pr(s^{N-1} \mathbf{s}_{\leftarrow}^R | \boldsymbol{\pi}_{\mathbf{x}^N}, x^N \mathbf{x}^R) \delta_{\Omega, \Omega_F} \delta_{\Psi, \Psi_F} \\ &= e^{\Omega} e^{\Psi} \sum_{s^0 \mathbf{s} \in \mathcal{S}^N} \Pr(s^{N-1} \mathbf{s}_{\leftarrow}^R | \boldsymbol{\pi}_{\mathbf{x}^N}, x^N \mathbf{x}^R) \delta_{\Omega, \Omega_F} \delta_{\Psi, \Psi_F} \\ &= e^{\Omega} e^{\Psi} \sum_{s^{N-1} \mathbf{s}_{\leftarrow}^R \in \mathcal{S}^N} \Pr(s^{N-1} \mathbf{s}_{\leftarrow}^R | \boldsymbol{\pi}_{\mathbf{x}^N}, x^N \mathbf{x}^R) \delta_{\Omega, -\Omega_R} \delta_{\Psi, -\Psi_R} \\ &= e^{\Omega} e^{\Psi} \Pr(-\Omega, -\Psi | \boldsymbol{\pi}_{\mathbf{x}^N}, x^N \mathbf{x}^R x^0) . \end{aligned}$$

Finally, we rewrite this to give the *extended CFT for irreversible processes*:

$$\frac{\Pr(\Omega, \Psi | \boldsymbol{\pi}_{\mathbf{x}^0}, x^0 \mathbf{x}^N)}{\Pr(-\Omega, -\Psi | \boldsymbol{\pi}_{\mathbf{x}^N}, x^N \mathbf{x}^R x^0)} = e^{\Psi} e^{\Omega} . \quad (23)$$

D. Interpretation

In the special case of isothermal time-symmetric driving— $x^0 \mathbf{x} x^0 = x^0 \mathbf{x}^R x^0 = x^0 x^1 x^2 \dots x^2 x^1 x^0$ —and starting from an equilibrium distribution, Eq. (23) provides a useful comparison between values of excess work achieved by the single time-symmetric driving protocol:

$$\frac{\Pr(W_{\text{ex}}, \Psi)}{\Pr(-W_{\text{ex}}, -\Psi)} = e^{\Psi} e^{\beta W_{\text{ex}}} . \quad (24)$$

Equation (23) should be compared to the original CFT that, in its most general form, can be written (with nec-

essary interpretation) as [35]:

$$\frac{\Pr_{\mathbf{F}}(\Omega)}{\Pr_{\mathbf{R}}(-\Omega)} = e^{\Omega} . \quad (25)$$

It is tempting to write Eq. (25) as:

$$\frac{\Pr(\Omega|\pi_{\mathbf{x}^0}, x^0 \mathbf{x} x^N)}{\Pr(-\Omega|\pi_{\mathbf{x}^N}, x^N \mathbf{x}^R x^0)} \stackrel{?}{=} e^{\Omega} . \quad (26)$$

This form presents some concerns, however. In the case of detailed balance, though, $\Psi = 0$ for all trajectories, and so our Eq. (23) guarantees Eq. (26) in the case of detailed balance. Crooks' original CFT derivation [2, 17] also assumed detailed balance, and so Eq. (26) was implied.

However, absent detailed balance, Eq. (25) has a rather different interpretation: $\Pr_{\mathbf{R}}(\cdot)$ then implies not only the reversed driving, but also that the distribution describes a different ‘‘reversed’’ system that is not of direct physical relevance [35, 36]. One consequence is that the probabilities in the numerator and denominator are not comparable in any physical sense. So, in general, we have:

$$\frac{\Pr(\Omega|\pi_{\mathbf{x}^0}, x^0 \mathbf{x} x^N)}{\Pr(-\Omega|\pi_{\mathbf{x}^N}, x^N \mathbf{x}^R x^0)} \neq e^{\Omega} . \quad (27)$$

In contrast, our irreversible CFT in Eq. (23) compares probabilities of entropy production (and path irreversibility) for the same thermodynamic system under a control protocol and under the reversed control protocol. Equation (23), unlike equalities involving an unphysical dual dynamic as in Eq. (25), allows a clear and meaningful physical interpretation of the relationship between entropies produced and, moreover, is not limited by assuming detailed balance.

Note that our Eq. (23), expressed in terms of excess environmental entropy production Ω and path irreversibility Ψ , does not make explicit mention of temperature. Indeed, if temperature dependence is folded into different environmental inputs x , then Eq. (23) applies just as well to systems driven by environments with spatially inhomogeneous temperature distributions that change in time. Explicitly, $\pi_{\mathbf{x}}$ and $\pi_{\mathbf{x}'}$ could represent the distribution over effective states induced by environmental conditions associated with x and x' including their different spatial distributions of temperature.

E. Translation to Steady-State Thermodynamics

A better understanding of the irreversible CFT comes by comparing it to recent related work. Most directly, our results complement those on driven transitions between

NESSs. Specifically, the importance of nondetailed-balanced dynamics in enabling the organization of complex nonequilibrium behavior has been considered previously. For example, Ref. [30] also introduced a path entropy which is an ensemble average of that considered here.

Another comparison is found in Ref. [14]'s nonequilibrium thermodynamics over NESSs using housekeeping Q_{hk} and excess Q_{ex} heats. While that treatment focused on Langevin dynamics, we find that in general Q_{hk} corresponds directly to our path irreversibility Ψ . Specifically, in the isothermal setting there, according to Eq. (35), we have:

$$\beta Q_{\text{hk}} \approx \Psi .$$

Indeed, for isothermal Markovian dynamics Eq. (7.7) of Ref. [37] suggests (via their Eqs. (2.11) and (7.1)) that this is in fact an equality:

$$\beta Q_{\text{hk}} = \Psi . \quad (28)$$

Reference [23] called the irreversibility Ψ the *adiabatic contribution* to entropy production. Several related translations from Ref. [14] to our setting can also be easily made: $\rho_{\text{ss}}(s; x) \rightarrow \pi_x(s)$, $\phi(s; x) \rightarrow -\log \pi_x(s)$, $\Delta S \rightarrow \Delta S_{\text{ss}}$, and $\beta Q_{\text{ex}} + \Delta \phi \rightarrow \Omega$. Hence, $\langle \Omega \rangle \geq 0$ (for Langevin systems) is Ref. [14]'s main result. From these connections, we see that our development not only provides new constraints on detailed fluctuations, but also extends these earlier results beyond Langevin systems.

Exposing these translations allows reformulating our detailed fluctuation theorems to steady-state thermodynamics (SST). We have:

$$e^{\gamma(s^0|\mu_{\mathbf{F}}, x^0) - \gamma(s^{N-1}|\mu_{\mathbf{R}}, x^N)} e^{\Omega_{\mathbf{F}} + \Psi_{\mathbf{F}}} = e^{\beta(Q_{\text{ex}} + Q_{\text{hk}}) + \Delta S^{\text{sys}}} = e^{\Delta S_{\mathbf{F}}^{\text{tot}}} ,$$

where $\Delta S^{\text{sys}} \equiv -\ln \frac{\Pr(s^{N-1}|\mu_{\mathbf{R}}, x^N)}{\Pr(s^0|\mu_{\mathbf{F}}, x^0)}$ when we choose $\langle \mu_{\mathbf{R}} | = \langle \mu_{\mathbf{F}} | \mathsf{T}(\mathcal{S} \rightarrow \mathcal{S} | x^0 \mathbf{x})$ and where $S_{\mathbf{F}}^{\text{tot}}$ is the total change in entropy in forward time. This yields:

$$\frac{\Pr(\mathcal{S}_{0:N} = s^0 \mathbf{s} | \mathcal{S}_{-1} \sim \mu_{\mathbf{F}}, X_{0:N} = x^0 \mathbf{x})}{\Pr(\mathcal{S}_{0:N} = s^{N-1} \mathbf{s}_{\mathbf{R}}^R | \mathcal{S}_{-1} \sim \mu_{\mathbf{R}}, X_{0:N} = x^N \mathbf{x}^R)} = e^{\Delta S_{\mathbf{F}}^{\text{tot}}} .$$

And so, we immediately see that:

$$\langle e^{-\Delta S_{\mathbf{F}}^{\text{tot}}} \rangle_{\Pr(\mathcal{S}_{0:N} = s^0 \mathbf{s} | \mathcal{S}_{-1} \sim \mu_{\mathbf{F}}, X_{0:N} = x^0 \mathbf{x})} = 1 .$$

This extends the validity of Ref. [38]'s general integral fluctuation theorem beyond Langevin dynamics. Since the total change in entropy is time asymmetric— $\Delta S_{\mathbf{F}}^{\text{tot}} = -\Delta S_{\mathbf{R}}^{\text{tot}}$ —we obtain the most direct CFT generalization

valid outside of detailed balance:

$$\frac{\Pr(\Delta S^{\text{tot}}|\boldsymbol{\pi}_{\mathbf{x}^0}, x^0 \mathbf{x}^N)}{\Pr(-\Delta S^{\text{tot}}|\boldsymbol{\pi}_{\mathbf{x}^N}, x^N \mathbf{x}^R x^0)} = e^{\Delta S^{\text{tot}}} . \quad (29)$$

Again, this does not invoke a dual, unphysical dynamic. Equation (29) has been reported previously in various settings; see, e.g., Eq. (21) of Ref. [23] and Eq. (43) of Ref. [36]. The result gives a detailed fluctuation relation for the change in total entropy production when transitioning between steady states.

The new detailed fluctuation theorem of Eq. (23) for joint distributions goes further in refining SST. If starting in a steady state and executing a protocol in an isothermal environment, we find that:

$$\frac{\Pr(W_{\text{ex}}, Q_{\text{hk}}|\boldsymbol{\pi}_{\mathbf{x}^0}, x^0 \mathbf{x}^N)}{\Pr(-W_{\text{ex}}, -Q_{\text{hk}}|\boldsymbol{\pi}_{\mathbf{x}^N}, x^N \mathbf{x}^R x^0)} = e^{\beta Q_{\text{hk}}} e^{\beta W_{\text{ex}}} .$$

This novel relation gives strong constraints on the thermodynamic behavior of systems driven between NESSs, since it constrains the *joint* distribution for excess work and housekeeping heat. Moreover, nonsteady-state additions to free energy are predicted when an experiment does not start in steady state.

In the special case of time-symmetric driving— $x^0 \mathbf{x} x^0 = x^0 \mathbf{x}^R x^0 = x^0 x^1 x^2 \dots x^2 x^1 x^0$ —and starting from an equilibrium distribution, the preceding expression reduces to a useful comparison between excess work values achieved by the single time-symmetric protocol:

$$\frac{\Pr(W_{\text{ex}}, Q_{\text{hk}})}{\Pr(-W_{\text{ex}}, -Q_{\text{hk}})} = e^{\beta Q_{\text{hk}}} e^{\beta W_{\text{ex}}} .$$

Similar results were recently derived in Ref. [39] under more restrictive assumptions for underdamped Langevin systems.

This all said, one must use caution and not always identify Ψ with βQ_{hk} . Most importantly, not all sources of irreversibility are naturally characterized as “heat”. Thinking of irreversibility on its own dynamical terms is best.

F. Integral fluctuation theorems

Integral fluctuation theorems in the absence of detailed balance, starting arbitrarily far from equilibrium, also follow straightforwardly. One generalization of the inte-

gral fluctuation theorem [40] is:

$$\begin{aligned} & \langle e^{-\beta W_{\text{diss}} - \Psi} \rangle_{\Pr(s_{0:N}|\boldsymbol{\mu}_{\text{F}}, \mathbf{x})} \\ &= \sum_{s_{0:N} \in \mathcal{S}^N} \Pr(s_{0:N}|\boldsymbol{\mu}_{\text{F}}, \mathbf{x}) e^{-\beta W_{\text{diss}} - \Psi} \\ &= \sum_{s_{0:N} \in \mathcal{S}^N} \Pr(s^0 \leftarrow^{x^1} \dots \leftarrow^{x^N} \boldsymbol{\mu}(\boldsymbol{\mu}_{\text{F}}, \mathbf{x}) | \boldsymbol{\mu}(\boldsymbol{\mu}_{\text{F}}, \mathbf{x}), \mathbf{x}^R) \\ &= 1 . \end{aligned} \quad (30)$$

If the input is stochastic, then averaging over the input also gives:

$$\langle e^{-\beta W_{\text{diss}} - \Psi} \rangle_{\Pr(x_{0:N}, s_{0:N}|\boldsymbol{\mu}_{\text{F}})} = 1 .$$

Note that this relation does *not* require the system to be in steady state at any time.

From the concavity of the exponential function, it is tempting to assert a corresponding generalized Second Law of SST as:

$$\langle W_{\text{diss}} \rangle \geq -\langle Q_{\text{hk}} \rangle . \quad (31)$$

Although Eq. (31) is true, notably it is neither a strong nor useful bound. Let’s address this. Note that:

$$\langle e^{-\Psi} \rangle_{\Pr(s_{0:N}|\boldsymbol{\mu}_{\text{F}}, \mathbf{x})} = 1 \quad (32)$$

and:

$$\langle e^{-\beta W_{\text{diss}}} \rangle_{\Pr(s_{0:N}|\boldsymbol{\mu}_{\text{F}}, \mathbf{x})} = 1 . \quad (33)$$

Both follow from the normalization of probabilities of the conjugate dynamic. Therefore, $\langle \Psi \rangle \geq 0$ and $\langle W_{\text{diss}} \rangle \geq 0$. And, hence $\langle Q_{\text{hk}} \rangle \geq 0$ as shown in Ref. [41]. So, Eq. (31) is devoid of utility. Nevertheless, Eq. (30) puts a novel constraint on the joint distributions of W_{diss} and Ψ .

Introducing an artificial conjugate dynamic following Ref. [35] and following the derivation there with ϕ/β in place of E , when starting in the steady state distribution $\boldsymbol{\pi}_{\mathbf{x}^0}$, we can show that:

$$\langle e^{-\Omega} \rangle_{\Pr(s_{0:N}|\boldsymbol{\pi}_{\mathbf{x}^0}, \mathbf{x})} = 1 , \quad (34)$$

which implies the restriction $\langle \Omega \rangle \geq 0$. Despite similar appearance, this result has meaning beyond Crooks’ derivation of the Jarzynski equality, as it now also applies *atop nonequilibrium steady states*. Recall that Ω has general meaning as in Eq. (6): $\Omega = \beta W_{\text{ex}} = -\beta Q_{\text{ex}} + \Delta\phi$. So, Eq. (34) becomes:

$$\langle e^{-\beta W_{\text{ex}}} \rangle_{\Pr(s_{0:N}|\boldsymbol{\pi}_{\mathbf{x}^0}, \mathbf{x})} = 1 . \quad (35)$$

Effectively, this is Ref. [14]’s relation that, with our sign

convention for Q_{ex} , implies:

$$\begin{aligned} \langle \Omega \rangle &= \langle -\beta Q_{\text{ex}} + \Delta \phi \rangle \\ &\geq 0, \end{aligned} \quad (36)$$

for processes that start in steady state.

However, using Eq. (33), we find a more precise constraint on expected excess entropy production whether or not the system starts in steady state:

$$\langle \Omega \rangle \geq \Delta \langle \gamma \rangle, \quad (37)$$

where the RHS can be positive or negative, but can only be negative if the system starts out of steady state. When starting in a steady state, this yields:

$$\begin{aligned} \langle \Omega \rangle &\geq \langle \gamma_{\text{final}} \rangle \\ &= D_{\text{KL}}[\text{Pr}(\mathcal{S}_t | \mathcal{S}_0 \sim \boldsymbol{\pi}_{\mathbf{x}_0}, x_{1:t+1}) || \boldsymbol{\pi}_{\mathbf{x}_t}], \end{aligned} \quad (38)$$

which is a stronger constraint than the previous result of Eq. (36), since the RHS is always positive for $\text{Pr}(\mathcal{S}_t | \mathcal{S}_0 \sim \boldsymbol{\pi}_{\mathbf{x}_0}, x_{1:t+1}) \neq \boldsymbol{\pi}_{\mathbf{x}_t}$. Eq. (39) extends the validity of the main result obtained in Ref. [42] to now include the possibility of starting in a nonequilibrium steady state and allowing for non-detailed-balanced dynamics. (Note that ‘ W_{diss} ’ in Ref. [42] corresponds to our W_{ex} —it is excess work that is not necessarily yet dissipated.)

Integral fluctuation theorems for systems with controlled or intrinsic feedback also directly follow, as we now show, extending the theory of feedback control to the setting of transitions between NESSs.

G. Fluctuation theorems with an auxiliary variable

Actions made by a complex thermodynamic system can couple back from the environment to influence the system’s future input. To achieve this, the system may express an auxiliary random variable Y_t —the current “output” that takes on the values $y \in \mathcal{Y}$ and is instantaneously energetically mute, but may influence the future input and so does have energetic relevance.

The variable Y_t could be measurement, output, or any other auxiliary variable that influences the state or input sequences. To be concise, we introduce a shorthand for the time-ordered sequences of random variables: $\vec{X} \equiv X_{0:N}$, $\vec{S} \equiv S_{0:N}$, and $\vec{Y} \equiv Y_{0:N}$. And, for particular realizations of the sequences: $\vec{x} \equiv x^{0:N}$, $\vec{s} \equiv s^{0:N}$, and $\vec{y} \equiv y^{0:N}$. When time reversing realizations, we let $\overleftarrow{x} = x^{N-1}x^{N-2} \dots x^1x^0$ and $\overleftarrow{s} = s^{N-1}s^{N-2} \dots s^1s^0$. To clarify further, \vec{x} appearing inside a probability implies $\vec{X} = \vec{x}$ and \overleftarrow{s} appearing inside a probability im-

plies $\overleftarrow{S} = \overleftarrow{s}$.

We quantify how much the auxiliary variable is independently informed from the state sequence—beyond what could be known if given only the initial distribution over states and the driving history—via the unaveraged conditional mutual information:

$$\begin{aligned} i(\overleftarrow{S}; \overleftarrow{Y} | \vec{x}, \boldsymbol{\mu}_F) &\equiv \ln \frac{\text{Pr}(\overleftarrow{S}, \overleftarrow{Y} | \vec{x}, \boldsymbol{\mu}_F)}{\text{Pr}(\overleftarrow{Y} | \vec{x}, \boldsymbol{\mu}_F) \text{Pr}(\overleftarrow{S} | \vec{x}, \boldsymbol{\mu}_F)} \\ &= \ln \frac{\text{Pr}(\overleftarrow{S}, \overleftarrow{Y}, \vec{x} | \boldsymbol{\mu}_F)}{\text{Pr}(\overleftarrow{Y}, \vec{x} | \boldsymbol{\mu}_F) \text{Pr}(\overleftarrow{S} | \vec{x}, \boldsymbol{\mu}_F)}. \end{aligned}$$

Note that averaging over the input, state, and auxiliary sequences gives the familiar conditional mutual information: $I[\overleftarrow{S}; \overleftarrow{Y} | \vec{X}, \boldsymbol{\mu}_F] = \langle i[\overleftarrow{S}; \overleftarrow{Y} | \vec{x}, \boldsymbol{\mu}_F] \rangle_{\text{Pr}(\vec{x}, \vec{s}, \vec{y} | \boldsymbol{\mu}_F)}$.

As detailed in the App. B:

$$e^{\beta W_{\text{diss}} + i[\overleftarrow{S}; \overleftarrow{Y} | \vec{x}, \boldsymbol{\mu}_F] + \Psi} = \frac{\text{Pr}(\overleftarrow{S}, \overleftarrow{Y}, \vec{x} | \boldsymbol{\mu}_F)}{\text{Pr}(\overleftarrow{Y}, \vec{x} | \boldsymbol{\mu}_F) \text{Pr}(\overleftarrow{S} | \vec{x}, \boldsymbol{\mu}_F)},$$

where $\boldsymbol{\mu}_R = \boldsymbol{\mu}(\boldsymbol{\mu}_F, \vec{x})$. This leads directly to the integral fluctuation theorem:

$$\left\langle e^{-\beta W_{\text{diss}} - i[\overleftarrow{S}; \overleftarrow{Y} | \vec{x}, \boldsymbol{\mu}_F] - \Psi} \right\rangle_{\text{Pr}(\vec{x}, \vec{y}, \vec{x} | \boldsymbol{\mu}_F)} = 1. \quad (40)$$

However, as before, the resulting bound on $\langle W_{\text{diss}} \rangle$ is not the tightest possible. Alternatively, we can invoke the normalization of conjugate dynamics to show:

$$\left\langle e^{-\beta W_{\text{diss}} - i[\overleftarrow{S}; \overleftarrow{Y} | \vec{x}, \boldsymbol{\mu}_F]} \right\rangle_{\text{Pr}(\vec{x}, \vec{y}, \vec{x} | \boldsymbol{\mu}_F)} = 1. \quad (41)$$

This implies a new lower bound for the revised Second Law of Thermodynamics:

$$\langle W_{\text{diss}} \rangle \geq -k_B T I[\overleftarrow{S}; \overleftarrow{Y} | \vec{X}, \boldsymbol{\mu}_F], \quad (42)$$

enabled by the conditional mutual information between state-sequence and auxiliary sequence, given input-sequence. Notably, this relation holds arbitrarily far from equilibrium and allows for the starting and ending distributions to be nonsteady-state.

We may also be interested in the unaveraged unconditioned mutual information between the auxiliary variable sequence and the joint input–state sequence. Then, using:

$$i(\overleftarrow{Y}; \overleftarrow{X} \overleftarrow{S} | \boldsymbol{\mu}_F) \equiv \ln \frac{\text{Pr}(\overleftarrow{Y}, \overleftarrow{S}, \overleftarrow{X} | \boldsymbol{\mu}_F)}{\text{Pr}(\overleftarrow{Y} | \boldsymbol{\mu}_F) \text{Pr}(\overleftarrow{X}, \overleftarrow{S} | \boldsymbol{\mu}_F)},$$

we find that, in general:

$$\langle W_{\text{diss}} \rangle \geq -k_B T I[\overleftarrow{Y}; \overleftarrow{X} \overleftarrow{S} | \boldsymbol{\mu}_F] \quad (43)$$

and when starting in steady-state:

$$\langle \Omega \rangle \geq -I[\vec{Y}; \vec{X}\vec{S} | \pi_{x^0}]. \quad (44)$$

One can now continue in this fashion to successively derive a seeming unending sequence of fluctuation theorems. Let's stop, though, with one more and discuss its interpretations and applications.

As a final set of example integral fluctuation theorems, we follow Ref. [3] in defining:

$$i_{\text{SU}} \equiv \ln \frac{\Pr(\vec{y}, \vec{s} | \mu_0)}{\Pr(\vec{y} | \mu_0) \Pr(\vec{s} | \mu_0, \vec{x})}.$$

(This is Ref. [3]'s I_C , if $\mu_0 \rightarrow \pi_{x^0}$.) Technically speaking, this is not a mutual information, even upon averaging. Then, we arrive at the integral fluctuation theorems:

$$\langle e^{-W_{\text{diss}} - i_{\text{SU}} - \Psi} \rangle_{\Pr(\vec{s}, \vec{y}, \vec{x} | \mu_0)} = 1$$

and

$$\langle e^{-W_{\text{diss}} - i_{\text{SU}}} \rangle_{\Pr(\vec{s}, \vec{y}, \vec{x} | \mu_0)} = 1.$$

When starting from a steady-state distribution, we have the most direct generalization of Ref. [3]'s feedback control result, but extended to not require detailed balance:

$$\langle e^{-\Omega - i_{\text{SU}}} \rangle_{\Pr(\vec{s}, \vec{y}, \vec{x} | \pi_{x^0})} = 1.$$

When starting from a NESS, this suggests that:

$$\langle \Omega \rangle \geq -I_{\text{SU}}. \quad (45)$$

When the dynamics are detailed balance, this naturally reduces to the well known results of Ref. [3] and others: $\langle W \rangle \geq \Delta F_{\text{eq}} - k_{\text{B}}T I_{\text{SU}}$.

In the feedback control setting, Y_n is said to be the random variable for measurements at time n . This suggests that Y_n is a function of \mathcal{S}_n and the outcome of Y_n effectively induces different Markov chains over the states since X_{n+1} is a function of Y_n —i.e., $x_{n+1}(y_n(s_n))$.

With our interest in complex autonomous systems, we note that our results give new bounds on the Second Law of Thermodynamics for highly structured complex systems strongly coupled to an environment. A preliminary application of this was presented out in Ref. [43]. We offer our own in the next section. Analysis of thermodynamic systems with the agency to influence their environmental input, via some kind of coupling or feedback, say, will likely benefit from our extended theory.

How can we reconcile this with other inequalities without auxiliary Y ? The other inequalities used averages of variable occurrence already conditioned on \vec{x} . However,

if input x and states s can influence each other dynamically through auxiliary y , then averaging over their joint dynamic allows less dissipation than the traditional Second Law suggests.

If \mathcal{S} represents the random variable for one subset of a system's degrees of freedom, and Y represents the random variable for another subset of a system's degrees of freedom, then the intrinsic nonextensivity of the thermodynamic entropy $S(\mathcal{S}, Y|X) = S(\mathcal{S}|X) + S(Y|X) - k_{\text{B}} I(\mathcal{S}; Y|X)$ goes a long way towards explaining the physics of information stimulating the recent resurgence of Maxwellian demonology. This viewpoint will be further developed elsewhere.

VI. NESS TRANSITIONS IN NEURONAL DYNAMICS

Acting in concert, voltage-gated sodium ion channels and potassium ion channels are the primary thermodynamic substrate that drives the evolution of membrane potentials in neurons [44]. Together, these voltage-gated channels are the primary generators of the action potentials or “spikes” that are the basic signals whose collective patterns support neural information processing [45]. In experiments, if the cell membrane is voltage clamped, then the channels approach a stationary distribution over their conformational states according to the effective energies of their biomolecular conformations at that voltage [46]. However, absent clamping, the channels influence their own voltage input dynamically through their current output. The result is spontaneous spiking patterns.

Although this is not the setting in which to analyze the full richness of ion channel interactions, we will use the sodium channel under different voltage-driving protocols as a relatively straightforward example to demonstrate the insights on NESS transitions gained from preceding theoretical results. That is, while potassium ion channels are somewhat structured, the sodium ion channel exhibits a more structured and so more illustrative dynamic over its coarse-grained state space of functional protein conformations.

A. Ion Channel Dynamics

For sodium ions to move through the neural membrane, a channel's activation gates must be open and the deactivation gate must not yet plug the channel [47]. The rates of transitions among the conformational states have a highly nontrivial dependence on voltage across the cell membrane. Beyond this voltage dependence, while

the activation gates act largely independently of one another, the inactivation gate cannot plug the channel until at least some of the activation gates are open. This causal architecture was not yet captured by the relatively macroscopic differential equations introduced in the pioneering work of Hodgkin and Huxley [46]. Since then, however, it has been summarized by experimentally-motivated voltage-dependent Markov chains over the causally relevant conformational states [48]. Here, we follow the model implied in Ref. [47], whose voltage-dependent Markov chain we show in Fig. 1.

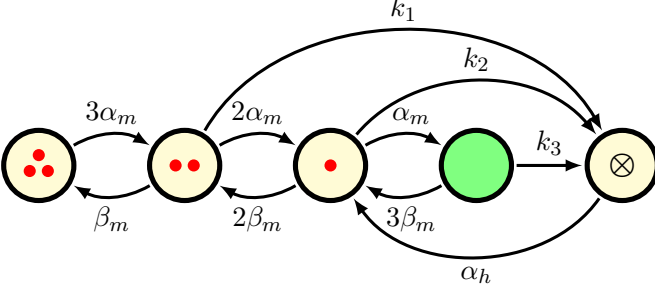


FIG. 1. Markov chain representation of the input-conditioned state-to-state rate matrix $G^{(\mathcal{S} \rightarrow \mathcal{S}|v)}$. Self-transitions are implied but, for tidiness, are not shown explicitly. These coarse-grainings of conformational states have biologically important functional interpretations. The number of (red) dots in each effective state corresponds to the number of activation gates that close off the channel. For example, when the channel is in the leftmost state, three activation gates are still active in blocking the channel. The solid (green) state corresponds to the channel being open. This is the only one of the five states in which sodium current can flow. The last state, marked \otimes , corresponds to channel inactivation by the inactivation gate—when the channel is plugged by its “ball and chain”. Subsequent figures use the state numbering 1 through 5 for state identification, which corresponds to enumeration of the states from left to right here.

As the mesoscopic system of thermodynamic interest, the voltage-dependent Markov chain can be reinterpreted as a transducer that takes in voltage $v \in \mathbb{R}$ across the cell membrane as its input and makes transitions over its conformational state space \mathcal{S} according to an infinite set of transition matrices $\{\mathbb{T}^{(\mathcal{S} \rightarrow \mathcal{S}|v)}\}_{v \in \mathbb{R}}$. Although the set is uncountable, the voltage-conditioned transition matrices are all described succinctly via time-independent functions of the voltage appearing in the transition elements; denoted α_m , β_m , and α_h . Time discretization of the continuous-time dynamic is straightforward and well behaved as $\Delta t \rightarrow 0$. If the voltage v is approximately constant during the infinitesimal interval Δt , then the state-to-state transition matrix is:

$$\begin{aligned} \mathbb{T}_{\Delta t}^{(\mathcal{S} \rightarrow \mathcal{S}|v)} &= e^{(\Delta t)G^{(\mathcal{S} \rightarrow \mathcal{S}|v)}} \\ &\approx I + (\Delta t)G^{(\mathcal{S} \rightarrow \mathcal{S}|v)}, \end{aligned}$$

where I is the identity matrix and $G^{(\mathcal{S} \rightarrow \mathcal{S}|v)}$ is the infinitesimal generator of time evolution:

$$G^{(\mathcal{S} \rightarrow \mathcal{S}|v)} \equiv \begin{bmatrix} -3\alpha_m & 3\alpha_m & 0 & 0 & 0 \\ \beta_m & -(2\alpha_m + \beta_m + k_1) & 2\alpha_m & 0 & k_1 \\ 0 & 2\beta_m & -(\alpha_m + 2\beta_m + k_2) & \alpha_m & k_2 \\ 0 & 0 & 3\beta_m & -(3\beta_m + k_3) & k_3 \\ 0 & 0 & \alpha_h & 0 & -\alpha_h \end{bmatrix}. \quad (46)$$

Specifically, α_m , β_m , and α_h are voltage-dependent variables, as found in the Hodgkin and Huxley model

[46, 47]:

$$\begin{aligned} \alpha_m(v) &= \frac{(v + 40 \text{ mV})/10 \text{ mV}}{1 - \exp[-(v + 40 \text{ mV})/10 \text{ mV}]}, \\ \beta_m(v) &= 4 \exp[-(v + 65 \text{ mV})/18 \text{ mV}], \end{aligned}$$

and:

$$\alpha_h(v) = \frac{7}{100} \exp[-(v + 65 \text{ mV})/20 \text{ mV}] .$$

See Fig. 2. The reaction-rate constants are $k_1 = 6/25 \text{ ms}^{-1}$, $k_2 = 2/5 \text{ ms}^{-1}$, and $k_3 = 3/2 \text{ ms}^{-1}$.

We developed new spectral decomposition methods from the meromorphic functional calculus [49, 50] to circumvent the inherent ill-conditioning in ion channel dynamics. Using these we can analytically calculate most, if not all, properties—e.g., dynamics, expected current, thermodynamics, information measures, and the like—about this model directly from the transition dynamic.

Since we are interested in thermodynamics, though, let us focus on determining the steady-state surprisals of the conformational states. For any persistent environmental input, the effective energies of the various conformational states are determined by their relative stationary occupation probability; according to Eq. (2), $\pi_x(s) = e^{-\phi(x,s)}$. The stationary distribution π_v induced by persistent v is the left eigenvector of $T_\tau^{(\mathcal{S} \rightarrow \mathcal{S}|v)}$ associated with the eigenvalue of unity. Equivalently, and more convenient in this case, π_v is the left eigenvector of $G^{(\mathcal{S} \rightarrow \mathcal{S}|v)}$ associated with the eigenvalue of zero. Via $\langle \pi_v | G^{(\mathcal{S} \rightarrow \mathcal{S}|v)} = \vec{0}$, we find that the steady state distribution for any persistent v is:

$$\begin{aligned} \pi_v \propto & \left(\frac{\beta_m}{3\alpha_m}, 1, \frac{2\alpha_m + k_1}{2\beta_m}, \frac{\alpha_m}{2\beta_m} \left(\frac{2\alpha_m + k_1}{3\beta_m + k_3} \right), \right. \\ & \left. \frac{1}{\alpha_h} \left[k_1 + \frac{2\alpha_m + k_1}{2\beta_m} \left(k_2 + \frac{k_3\alpha_m}{3\beta_m + k_3} \right) \right] \right) \\ & \propto e^{-\phi(v)} , \end{aligned}$$

which immediately yields the steady-state surprisals for conformational states at a constant environmental input v . The steady-state surprisal is shown for each conformational state in Fig. 3, as a function of the voltage-clamped membrane potential v .

B. Ion channel (ir)reversibility

Recall that detailed balance is the condition that, for states a and b and environmental input x , $\Pr(b \xrightarrow{x} a) / \Pr(a \xrightarrow{x} b) = \pi_x(a) / \pi_x(b)$. Interestingly, in this biologically inspired model, detailed balance is satisfied by *several, but not all* of the state-transition pairs.

For example, for very small Δt :

$$\begin{aligned} \frac{\Pr(1 \xrightarrow{v} 2)}{\Pr(2 \xrightarrow{v} 1)} &= \frac{3\alpha_m}{\beta_m} \\ &= \frac{\pi_v(2)}{\pi_v(1)} . \end{aligned}$$

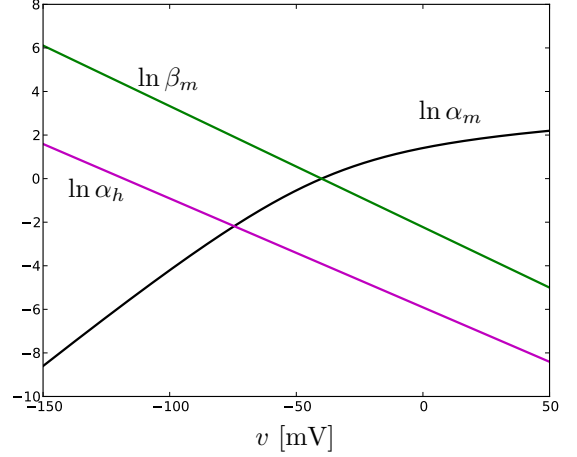


FIG. 2. Markov transition parameter voltage dependencies: Plots of $\ln \alpha_m$, $\ln \beta_m$, and $\ln \alpha_h$. These plots show that at -100 mV , $\beta_m \gg \alpha_h \gg \alpha_m \approx 0$. At $+10 \text{ mV}$, $\alpha_m \gg \beta_m \gg \alpha_h \approx 0$.

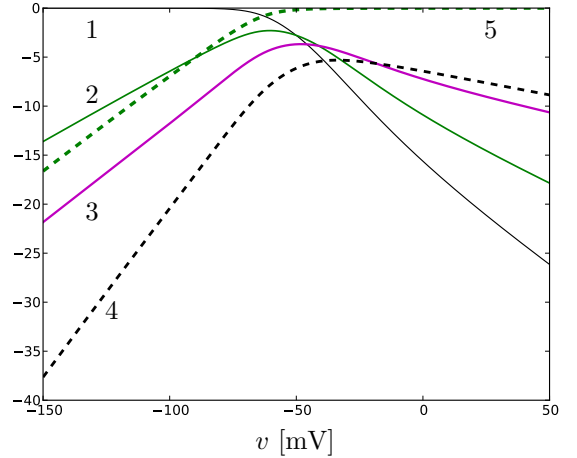


FIG. 3. Steady-state distribution voltage dependence: Negative of the steady-state surprisals, $\ln \pi_v(s) = -\phi(v, s)$ for each conformational state. Each curve labeled by the state-number to which it corresponds. Note that -100 mV and $+10 \text{ mV}$ (relevant for later) are extremes in that $\pi_{v_a} \approx \delta_1$ and $\pi_{v_b} \approx \delta_5$.

That is, this transition pair satisfies detailed balance. Hence, all transitions between these states are completely reversible: $\Psi(2|1, v) = \Psi(1|2, v) = \Psi(22121112|2, v) = 0$.

However, this does not hold for other transition pairs.

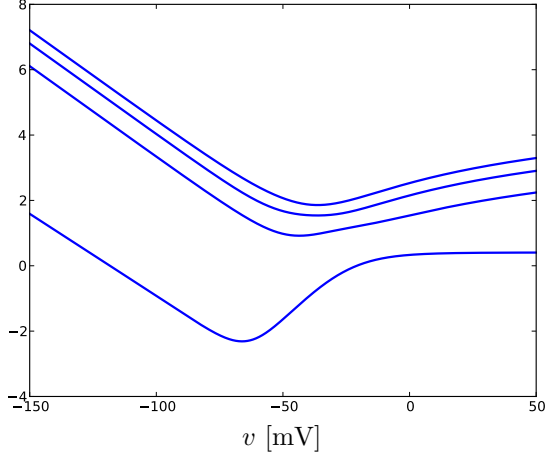


FIG. 4. Modes of the state-to-state dynamic: All eigenvalues of G are real and nonpositive. There is a zero eigenvalue associated with stationarity and four negative eigenvalues associated with decay rates from the states. Plots of $\ln[-\lambda(v)]$ for $\lambda(v) \in \Lambda_{G(\mathcal{S} \rightarrow \mathcal{S}|v)}$. Smaller $\ln(-\lambda)$ corresponds to longer time-scales. The zero eigenvalue maps to $-\infty$.

Consider transitions between states 2 and 3:

$$\begin{aligned} \frac{\Pr(2 \xrightarrow{v} 3)}{\Pr(3 \xrightarrow{v} 2)} &= \frac{\alpha_m}{\beta_m} \\ &\neq \frac{\pi_v(3)}{\pi_v(2)} \\ &= \frac{\alpha_m + k_1/2}{\beta_m} . \end{aligned}$$

Most other transitions also violate detailed balance.

Since detailed balance does not hold for the effective dynamic, the theory developed above is essential to analyzing the sodium ion channel thermodynamics. Moreover, the fact that $G(\mathcal{S} \rightarrow \mathcal{S}|v)$ has null entries that are nonzero for its transpose implies that paths involving these transitions will be infinitely irreversible: $\Psi = \infty$ for such paths as $\Delta t \rightarrow 0$; specifically, the transitions of $G(\mathcal{S} \rightarrow \mathcal{S}|v)$ with rates k_1 and k_3 . Forbidden transitions are an extreme form of irreversibility that are nevertheless commonly observed for complex systems, as the ion channel so readily illustrates. In it, the asymmetry in allowed transitions can be traced to different *mechanisms* facilitating different paths through the state space. Whether the irreversibility is truly infinite or just practically infinite does not matter much for the excess thermodynamics, although it will of course affect the calculated distribution of Ψ . (Conventional, linear algebraic methods are inadequate to overcome these technical challenges. The spectral decomposition methods, mentioned above, are required.)

C. Step Function Drive

With this understanding of ion channel NESSs, let's now turn to the thermodynamics induced by driving between them. We first consider the particular voltage protocol of $v_a \equiv -100$ mV for all time except a $v_b \equiv 10$ mV pulse for 5 ms starting at $t = 0$. This is an example of continuous-time dynamics and deterministic driving. The system begins equilibrated with the static environmental drive $v_a = -100$ mV. The initial distribution over \mathcal{S} is thus $\boldsymbol{\mu}_0 = \boldsymbol{\pi}_{v_a}$, where $\boldsymbol{\pi}_{v_a}$ is the left eigenvector of $G(\mathcal{S} \rightarrow \mathcal{S}|v_a)$ associated with the eigenvalue of zero.

During an epoch of fixed $v = V$, the net transition dynamic after τ ms becomes:

$$\mathcal{T}_\tau(\mathcal{S} \rightarrow \mathcal{S}|v=V) = e^{\tau G(\mathcal{S} \rightarrow \mathcal{S}|v=V)} .$$

Therefore, the distribution over states induced by the driving protocol is:

$$\langle \boldsymbol{\mu}_t | = \begin{cases} \langle \boldsymbol{\pi}_{v_a} | & \text{for } t \leq 0 \\ \langle \boldsymbol{\pi}_{v_a} | e^{t G_b} & \text{for } 0 < t \leq 5 \text{ ms} \\ \langle \boldsymbol{\pi}_{v_a} | e^{5 G_b} e^{(t-5) G_a} & \text{for } t > 5 \text{ ms} \end{cases} , \quad (47)$$

where, for brevity, we defined: $G_a \equiv G(\mathcal{S} \rightarrow \mathcal{S}|v_a)$ and $G_b \equiv G(\mathcal{S} \rightarrow \mathcal{S}|v_b)$. These are especially useful when expressing the rate matrix via its spectral decomposition, using the methods of Refs. [49, 50]. Besides the zero eigenvalue, there are only four other eigenvalues of G that are determined via $\det(\lambda I - G) = 0$.

Figure 4 shows G 's eigenvalues as a function of v , which indicates the voltage-dependent timescales of probability decay for modes of occupation probability. The associated decay rates play a prominent role in Fig. 5, which shows the time-dependent distribution induced over states by the 5 ms voltage pulse driving protocol.

Having the distribution over states at all times is powerful knowledge. For example, since the current through a single channel is binary—either 0 or $I_0(v)$ —and since current only flows in the open conformation, the expected current through the channel is $I_0(v) = g_0[v - V_{\text{Na}}]$ times the expectation value of being in the open state:

$$\langle I(t) \rangle = g_0 [v(t) - V_{\text{Na}}] \langle \boldsymbol{\mu}(t) | \delta_{\text{open}} \rangle ,$$

where $\delta_{\text{open}} = (0, 0, 0, 1, 0)$, g_0 is the conductance of an open Na^+ channel, and $V_{\text{Na}} = \frac{k_B T}{e^+} \ln \frac{[\text{Na}^+]_{\text{out}}}{[\text{Na}^+]_{\text{in}}} \approx 90$ mV is the Nernst potential for sodium in a typical mammalian neuron [47]. To be clear $\langle I(t) \rangle$ is what would be observed from an ensemble of channels in a local patch of cell membrane experiencing the same driving. The current produced from the Markovian model appears to be more realistic than what would be expected

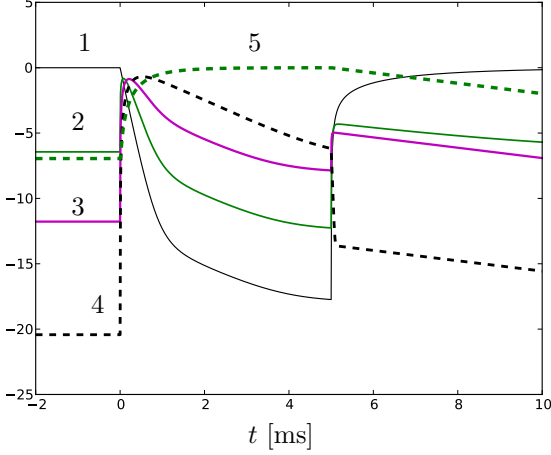


FIG. 5. Na^+ ion channel NESS transitions: temporal evolution of the distribution μ_t of ion-channel conformational states induced by a deterministic 5 ms voltage pulse is shown via plots of $\ln \langle \mu_t | s \rangle$ for all $s \in \mathcal{S}$. Curves are labeled by the state-number to which each component corresponds.

from the Hodgkin–Huxley model [47]. Moreover, using our spectral-decomposition methods for functions of a Markov chain [49, 50], this current can now be obtained in closed-form.

Let us start the thermodynamic investigation by considering excess work W_{ex} . With $\tau = N\Delta t$, we take the limit of $\Delta t \rightarrow 0$ while keeping the product $N\Delta t = \tau$ constant. Then the expected excess work per $k_B T$, from time t_0 to time $t_0 + \tau$, is:

$$\beta \langle W_{\text{ex}} \rangle = \int_{t_0}^{t_0 + \tau} \langle \mu_t | d\phi_{v(t)} / dt \rangle dt .$$

However, it should be clear that, for this stepped voltage protocol, excess work is *only performed on this system at the very onset and subsequently at the end* of the step driving. Indeed, this is the only time that the driving $v(t)$ changes and, thus, the only time that the state-dependent rate of work $d|\phi_{v(t)}|/dt$ is nonzero. As we let $\Delta t \rightarrow 0$, the expected excess work (divided by $k_B T$) near the onset of driving becomes a step function with height:

$$\begin{aligned} \lim_{\epsilon \rightarrow 0^+} \langle \Omega(t = \epsilon) - \Omega(t = -\epsilon) \rangle \\ = \sum_{s \in \mathcal{S}} \langle \pi_{v_a} | s \rangle [\phi(10\text{mV}, s) - \phi(-100\text{mV}, s)] , \end{aligned}$$

where $\langle \pi_{v_a} | s \rangle = \pi_{-100 \text{ mV}}(s)$.

Indeed, for this singular event, the full distribution of work performed can be given according to the probabilities that the system was in a particular state when the driving was applied. For $0 < t < 5$ ms, the probability

density function for βW_{ex} is:

$$p(\Omega) = \sum_{s \in \mathcal{S}} \pi_{-100 \text{ mV}}(s) \times \delta\left(\Omega - [\phi(10\text{mV}, s) - \phi(-100\text{mV}, s)]\right) ,$$

where $\delta(\cdot)$ here is the Dirac delta function. For $t > 5$ ms, the full excess environmental entropy production (EEEP) probability density function (pdf) is:

$$p(\Omega) = \sum_{s, s' \in \mathcal{S}} \langle \pi_{v_a} | s \rangle \langle s | e^{5G_b} | s' \rangle \times \delta\left(\Omega - [\phi(v_b, s) - \phi(v_a, s)] - [\phi(v_a, s') - \phi(v_b, s')]\right) .$$

From the Dirac delta function's argument and the sum over s and s' , it is clear that every nonzero-probability EEEP value Ω also has a nonzero probability for the negative $-\Omega$ of that EEEP value.

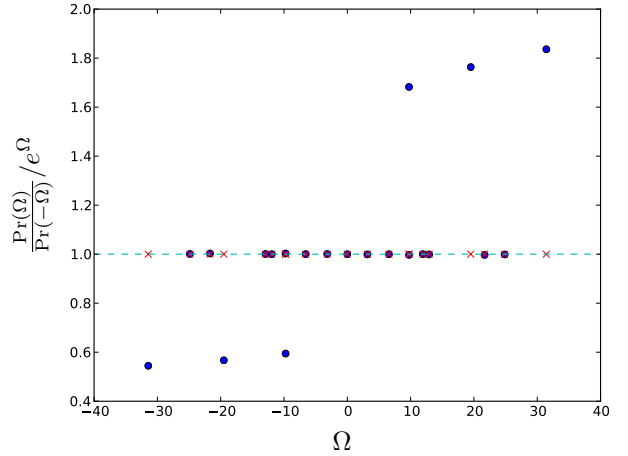


FIG. 6. Deviations from and agreement with the Crooks Fluctuation Theorem for nondetailed balance dynamics: Exact calculation of $\frac{\text{Pr}(\Omega)}{\text{Pr}(-\Omega)}/e^\Omega$ (blue dots) for all allowed values of $\Omega = \beta W_{\text{ex}}$ during the pulse drive. Since the system starts in equilibrium and since the driving is time-symmetric, a naive CFT interpretation suggests that all values lie at unity, dashed line (blue) and marked with a \times (red) wherever an allowed excess work value appears. Interestingly, many of the allowed work values *do* still fall on or very near unity. Absent detailed balance, though, Eq. (23) must be used to account for the actual distribution of excess environmental entropy production and path irreversibilities that, in addition to all other values, yields the six deviant markings (blue dots) above and below unity.

For the time-symmetric 5 ms voltage-pulse driving, Eq. (24) tells us that $\text{Pr}(\Omega, \Psi)/\text{Pr}(-\Omega, -\Psi) = e^\Psi e^\Omega$. Since there are infinitely many Ψ values to account for, we do not plot the joint distribution explicitly. However, we can appreciate the necessity of the relationship

by comparing it to the naive CFT interpretation that, for this case, suggests $\Pr(\Omega)/\Pr(-\Omega) = e^\Omega$. Figure 6 compares these by plotting $e^{-\Omega} \Pr(\Omega)/\Pr(-\Omega)$.

Allowed values of the excess work that do *not* lie on $e^{-\Omega} \Pr(\Omega)/\Pr(-\Omega) = 1$ demonstrate deviations from the naive CFT interpretation. Since the constant-voltage steady states are nonequilibrium and, thus, not microscopically reversible—i.e., $\Psi \neq 0$ for some state paths—the naive CFT interpretation cannot be true despite the time-symmetric driving. Perhaps the most surprising feature in Fig. 6 is that many of the probability ratios still *do* (almost) fall on the naive CFT line at unity. In part, this is due to a subset of the cycles in the NESS dynamic obeying detailed balance. Another contributing factor is that longer durations τ of fixed v induces a *net* dynamic $e^{\tau G}$ that *approaches* a detailed-balanced dynamic. That the values in Fig. 6 are sensible can be verified by checking the ratio of the joint probabilities $\langle \pi_{v_a} | s \rangle \langle s | e^{5G_b} | s' \rangle$ to the value of the joint probability with s and s' swapped.

In stark contrast to the instantaneous work contribution, the system's excess heat Q_{ex} unfolds over time, exhibiting a rich structure governed by the trajectories through the conformational state-space. The expected excess heat per $k_B T$ is:

$$\beta \langle Q_{\text{ex}} \rangle = \int_{t_0}^{t_0+\tau} \langle \dot{\mu} | \phi_{v(t)} \rangle dt . \quad (48)$$

over a duration τ , if starting at time t_0 . Since $v(t)$ is constant except at the two instants of change, the integral is easily solved exactly using the fundamental theorem of calculus and Eq. (47).

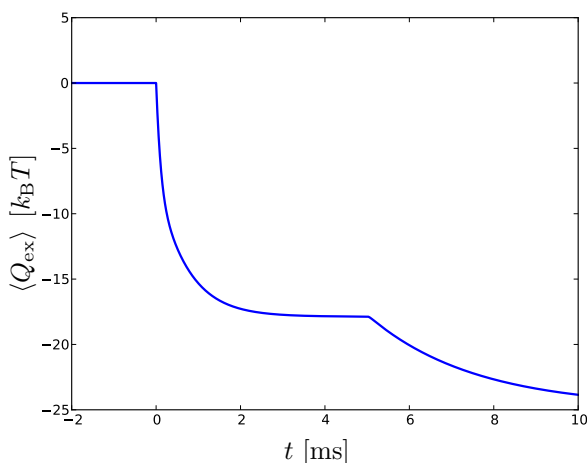


FIG. 7. Excess heat Q_{ex} per ion channel for an ensemble of ion channels embedded in a local patch of cell membrane; in units of $k_B T \approx 26$ meV. The two bouts of relaxation correspond to the ion channel adapting to sudden changes in voltage across the cell membrane.

For $t_0 < t_0 + \tau < 0$, the system is in the initial steady-state and has a time-independent heat and so a constant excess heat rate that vanishes:

$$\frac{\langle Q_{\text{ex}} \rangle}{\tau} = 0 .$$

Figure 7 shows the expected excess heat $\langle Q_{\text{ex}} \rangle$ over the course of the voltage-drive protocol. The steady-state average rate of excess heat production within any steady state is necessarily zero. However, the channel macromolecule responds to changes in the environment via conformational changes and corresponding heat productions that unfold on the timescale of milliseconds [51]. Notably, the expected heat is on the order of tens of $k_B T$, which for mammalian neurons is $k_B T = 1/\beta \approx 26$ meV.

For $0 = t_0 < t_0 + \tau < 5$ ms, the system has a time-varying heat production as it synchronizes to the new nonequilibrium steady state. From Eq. (47), we know that $\langle \mu_t \rangle = \langle \pi_{v_a} | e^{tG_b}$ in this case. At the same time, the steady-state-surprisal $|\phi_{v(t)}\rangle$ is time-independent during this epoch since v is temporarily fixed at $v_b = 10$ mV.

Again, as seen in Fig. 7, the expected excess heat drops for several milliseconds after the final voltage switch at 5 ms, as the ion channel re-adapts to its original steady state. For this second bout of relaxation the adaptation is slower since, in accordance with Fig. 4, the slowest timescale at $v_a = -100$ mV is slower than the slowest timescale at $v_b = 10$ mV.

Overall, we see that the excess thermodynamic quantities are well behaved and accessible: without needing to know the background biological upkeep of the Na^+ ion channel, we can access and control coarse degrees of freedom of the channel macromolecule via modulating the voltage across the cell membrane. Moreover, for the Na^+ channel, state-measurement and feedback on the timescale of milliseconds would allow significant alterations of heat and entropy production. Fortuitously, this suggests an accessible platform for laboratory experimentation. Next, we comment briefly on an intrinsic type of measurement and feedback that happens in vivo every moment.

D. Intrinsic feedback

Having come this far, we close illustrating the thermodynamics of NESS transitions with a final application. In a biologically active (in vivo) neuron, the input membrane voltage at each time depends on integrated current—a functional of the state distribution—up to that time. Our relations for modified integral fluctuation relations describe the thermodynamic agency of Na^+ channels in vivo, whereas conventional fluctuation rela-

tions fall short. Although there is certainly feedback in vivo, it is not the “feedback control” discussed recently. Importantly, no “outsider” forces the feedback; the feedback is intrinsic—woven into the system–environment joint dynamic. We leave a thorough investigation of the thermodynamics of intrinsic feedback to elsewhere. The success here, however, already suggests investigating other natural systems with intrinsic feedback—in joint nonlinear dynamics and complex networks—to test the new fluctuation theorems and computational methods in a broader class of interacting complex nonequilibrium systems.

VII. DISCUSSION

In light of our refined detailed fluctuation theorem Eq. (23) for nondetailed-balanced dynamics, we referred to the common belief in Eq. (27) as the naive CFT interpretation since it appeals to a nonphysical conjugate dynamics, as in Eq. (25). Similarly, we referred to failures of Eq. (26) as CFT violations. Nonetheless, with proper interpretation using the unphysical conjugate dynamics, Eq. (25) is mathematically correct even without detailed balance and can be a useful device for establishing integral FTs.

The CFT is often misinterpreted, though, despite receiving widespread attention. For example, Ref. [18] appealed to the naive CFT interpretation in the case of nondetailed balance dynamics of self-replication. However, as we showed, such applications are not valid. So, the statistical physics of self-replication either depends on an assumption of detailed balance or deserves a generalization. We believe the latter should be straightforward using our results, new spectral methods, and a derivation paralleling Ref. [18].

We hope that our nonintegral FTs—especially Eq. (23) that constrains the joint distribution of excess and house-keeping entropies—will provide better physical intuition for the structure of effective dynamics outside detailed balance. Path irreversibility clearly plays a prominent role. Although the preceding introduced a unifying framework, certain subclasses of path irreversibility have already been proposed recently.

In certain applications, for example, path irreversibility is governed by differences in chemical potential. In such cases, the irreversibility is quantitatively related to cycle affinities; see, for example, Ref. [52]. It essentially discovered a special case of the results developed here specifically applicable to the interesting example of a kinesin motor protein.

Recently, Ref. [53] elaborated on one type of irreversibility, called *absolute irreversibility*, that at first ap-

pears to constitute an extreme contribution to the total path irreversibility Ψ . This indeed is one interpretation, but not the full story. On closer examination, its result appears to coincide most directly with Eq. (33) which must be used when starting in a nonsteady-state, rather than with a violation of Eq. (34) which is simply inapplicable when starting in a nonsteady-state. We reinterpret that work as a testament to the importance of the nonsteady-state contribution to free energy change, $\beta^{-1}\Delta(\gamma)$. Explicitly:

$$\begin{aligned} \langle e^{-\beta W_{\text{diss}}} \rangle_{\text{Pr}(s_{0:N}|\boldsymbol{\mu}_F, \mathbf{x})} &= 1 \\ &\rightarrow \langle \Omega \rangle \geq \Delta \langle \gamma \rangle . \end{aligned}$$

This captures, for example, the entropy change associated with free expansion. From our viewpoint, however, any absolute irreversibility is only one extreme of the broader generalization introduced above to explore the consequences of irreversibility and nonsteady-state additions to free energy.

To frame our results in yet another way, we note that the “feedback control” imposed by an experimenter on an otherwise detailed-balanced system is a rather limited form of CFT violation. Yet it appears to be the only one having gained much recognition. This is odd. Hysteresis, to take one example, common in paradigmatically complex physical systems, provides a more physical manifestation of CFT violation. Even this is still a relatively tame deviation from detailed balance. Living systems are the true flagship of complex physical agents with intrinsic computational feedback across many levels of their organization. Our fluctuation relations describe all of these aspects, together.

In particular, they suggest how a system’s intrinsic model of its environment, together with an action policy that leverages knowledge captured in the model to control the environment, allows the system to play the survival game to its thermodynamic advantage. For example, an agent can use information about the environment to increase its nonsteady-state free energy and perform useful work—a phenomenon that is not only reminiscent of living beings, but also comes very near to defining them.

We hope that our results and methods stimulate investigating the excess thermodynamics of systems with intrinsic feedback—from designed “toy demons” to complex biological molecules affected by and simultaneously affecting their environments. Several biological examples that suggest themselves include kinesin motors [54], drug-operated channels [55], and dynamic synapses [56], just to name a few.

VIII. CONCLUSION

We presented generalized fluctuation theorems for complex nonequilibrium systems driven between NESSs. In addition to the detailed FTs that constrain joint distributions of excess and housekeeping quantities, we introduced integral fluctuation theorems in the presence of an auxiliary variable. The auxiliary variable need not be measurement nor any other meddling of an outsider. Due to this, it generalizes the theory of “feedback control” to the setting of arbitrary intrinsic feedback between system and environment.

A sequel to the above derives exact closed-form expressions for the moments of excess heat and excess work when the joint system–environment dynamic is governed by a (finite or countably infinite) discrete- or continuous-time hidden Markov model. A joint system can always be modeled as a joint hidden Markov model—at least as an approximation to the true joint dynamics. For this reason, our exact results should provide broadly applicable tools. The latter have particular theoretical advantage in giving access to what occurs in transient and asymptotic dynamics of excess thermodynamic quantities atop NESSs.

In summary, the traditional laws of thermodynamics are largely preserved for the renormalized “excess” thermodynamic quantities that arise naturally when considering nondetailed-balanced complex systems. However, the laws must be modified by the entropic contribution of path irreversibility. We noted that the latter turns out to be equivalent to steady-state thermodynamics’ housekeeping entropy.

Our relations still hold for excursions between equilibrium steady states, but we then have the simplification that $\Psi = \beta Q_{\text{hk}} = 0$. Consistently, equilibrium thermodynamics is a reduction of the theory of excess thermodynamic quantities with no housekeeping terms—when all paths are microscopically reversible.

Layers of emergence, typical of the biological world [57, Fig. 6], beg renormalization in terms of a hierarchy of housekeeping backgrounds [58]. The opportunity offered up by emergent levels of novel organization is a new richness in nondetailed-balanced effective dynamics—dynamics and structure that can be exploited by intelligent thermodynamic agency [59, 60]. We consider the thermodynamics of agency in a sequel, analyzing a simple autonomous agent that harvests energy by leveraging hidden correlations in a fluctuating environment.

We leave the development for now, but with an encouraging lesson: Even in nonstationary nonequilibrium, there is excess thermodynamic structure at any level of observation that we can access, control, and harness.

ACKNOWLEDGMENTS

We thank Tony Bell, Alec Boyd, Gavin Crooks, Chris Jarzynski, John Mahoney, Dibyendu Mandal, and Adam Rupe for useful feedback. We thank the Santa Fe Institute for its hospitality during visits. JPC is an SFI External Faculty member. This material is based upon work supported by, or in part by, the U. S. Army Research Laboratory and the U. S. Army Research Office under contracts W911NF-12-1-0234, W911NF-13-1-0390, and W911NF-13-1-0340.

Appendix A: Extension to Non-Markovian Instantaneous Dynamics

Commonly, theoretical developments assume state-to-state transitions are instantaneously Markovian given the input. This assumption works well for many cases, but fails in others with strong coupling between system and environment. Fortunately, we can straightforwardly generalize the results of stochastic thermodynamics by considering a system’s observable states to be functions of latent variables \mathcal{R} . The goal in the following is to highlight the necessary changes, so that it should be relatively direct to adapt our derivations to the non-Markovian dynamics.

a. Latent states, system states, and their many distributions

Even with constant environmental input, the dynamic over a system’s states need not obey detailed balance nor exhibit any finite Markov order. We assume that the classical observed states \mathcal{S} are functions $f : \mathcal{R} \rightarrow \mathcal{S}$ of a latent Markov chain. We also assume that the stochastic transitions among latent states are determined by the current environmental input $x \in \mathcal{X}$, which can depend arbitrarily on all previous input and system-state history. The Perron–Frobenius theorem guarantees that there is a stationary distribution over latent states associated with each fixed input x ; the function of the Markov chain maps this stationary distribution over latent states into the stationary distribution over system states. These are the stationary distributions associated with system NESSs.

We assume too that the \mathcal{R} -to- \mathcal{R} transitions are Markovian given the input. However, different inputs induce different Markov chains over the latent states. This can be described by a (possibly infinite) set of input-conditioned transition matrices over the latent state set \mathcal{R} : $\{\mathbb{T}(\mathcal{R} \rightarrow \mathcal{R}|x)\}_{x \in \mathcal{X}}$, where $\mathbb{T}_{i,j}^{(\mathcal{R} \rightarrow \mathcal{R}|x)} = \Pr(\mathcal{R}_t = r^j | \mathcal{R}_{t-1} = r^i, X_t = x)$. Probabilities regarding actual

state paths can be obtained from the latent-state-to-state transition dynamic together with the observable-state projectors, which we now define.

We denote distributions over the latent states as bold Greek symbols, such as $\boldsymbol{\mu}$. As in the main text, it is convenient to cast $\boldsymbol{\mu}$ as a row-vector, in which case it appears as the bra $\langle \boldsymbol{\mu} |$. The distribution over latent states \mathcal{R} implies a distinct distribution over observable states \mathcal{S} . A sequence of driving inputs updates the distribution: $\boldsymbol{\mu}_{t+n}(\boldsymbol{\mu}_t, x_{t:t+n})$. In particular:

$$\begin{aligned} \langle \boldsymbol{\mu}_{t+n} | &= \langle \boldsymbol{\mu}_t | \mathsf{T}(\mathcal{R} \rightarrow \mathcal{R} | x_{t:t+n}) \\ &= \langle \boldsymbol{\mu}_t | \mathsf{T}(\mathcal{R} \rightarrow \mathcal{R} | x_t) \mathsf{T}(\mathcal{R} \rightarrow \mathcal{R} | x_{t+1}) \dots \mathsf{T}(\mathcal{R} \rightarrow \mathcal{R} | x_{t+n-1}) . \end{aligned}$$

(Recall that time indexing is denoted by subscript ranges $n : m$ that are left-inclusive and right-exclusive.) An infinite driving history \vec{x} induces a distribution $\boldsymbol{\mu}(\vec{x})$ over the state space, and $\boldsymbol{\pi}_x$ is the specific distribution induced by tireless repetition of the single environmental drive x . This is the so-called ‘‘equilibrium distribution’’ associated with equilibrating with the environmental drive x . Explicitly:

$$\langle \boldsymbol{\pi}_x | = \lim_{n \rightarrow \infty} \langle \boldsymbol{\mu}_0 | \left(\mathsf{T}(\mathcal{R} \rightarrow \mathcal{R} | x) \right)^n .$$

Usefully, $\boldsymbol{\pi}_x$ can also be found as the left eigenvector of $\mathsf{T}(\mathcal{R} \rightarrow \mathcal{R} | x)$ associated with the eigenvalue of unity:

$$\langle \boldsymbol{\pi}_x | = \langle \boldsymbol{\pi}_x | \mathsf{T}(\mathcal{R} \rightarrow \mathcal{R} | x) . \quad (\text{A1})$$

The canonical equilibrium probabilities are this vector’s projection onto observable states: $\pi_x(s) = \langle \boldsymbol{\pi}_x | s \rangle$, where $|s\rangle = |\delta_{s,f(r)}\rangle$ has a vector-representation in the latent-state basis with elements of all 0s except 1s where the latent state maps to the observable state s .

Assuming latent-state-to-state transitions are Markovian allows the distribution $\boldsymbol{\mu}$ over these latent states to summarize the causal relevance of the entire driving history.

b. Implications

A semi-infinite history induces a particular distribution over system latent states and implies another particular distribution over its observable states. This can be usefully recast in terms of the ‘‘start’’ (or initial) distribution $\boldsymbol{\mu}_0$ induced by the path $x_{-\infty:1}$ and the driving history $x_{1:t+1}$ since then, giving the entropy of the induced state distribution:

$$\begin{aligned} h^{(s|\boldsymbol{\mu}_0, x_{1:t+1})} &= -\ln \Pr(\mathcal{S}_t = s | \boldsymbol{\mu}_0, x_{1:t+1}) \\ &= -\ln \langle \boldsymbol{\mu}_0 | \mathsf{T}(\mathcal{R} \rightarrow \mathcal{R} | x_{1:t+1}) | s \rangle . \end{aligned}$$

Or, employing the new distribution and the driving history since then, the path entropy (functional of state and driving history) can be expressed simply in terms of the *current* distribution over latent states and the candidate observable state s :

$$\begin{aligned} h^{(s|\boldsymbol{\mu})} &= -\ln \Pr(\mathcal{S}_t = s | \mathcal{R}_t \sim \boldsymbol{\mu}) \\ &= -\ln \langle \boldsymbol{\mu} | s \rangle . \end{aligned}$$

Averaging the path-conditional state entropy over observable states again gives a genuine input-conditioned Shannon state entropy:

$$\langle h^{(s_t|\vec{x}_t)} \rangle_{\Pr(s_t|\vec{x}_t)} = \mathsf{H}[\mathcal{S}_t | \overleftarrow{X}_t = \overleftarrow{x}_t] .$$

It is again easy to show that the state-averaged path entropy $k_B \mathsf{H}[\mathcal{S}_t | \overleftarrow{x}_t]$ is an extension of the system’s steady-state nonequilibrium entropy. In steady-state, the state-averaged path entropy reduces to:

$$\begin{aligned} k_B \mathsf{H}[\mathcal{S}_t | \overleftarrow{X}_t = \dots xxx] &= -k_B \mathsf{H}[\mathcal{S}_t | \mathcal{R}_t \sim \boldsymbol{\pi}_x] \\ &= -k_B \sum_{s \in \mathcal{S}} \pi_x(s) \ln \pi_x(s) \\ &= S_{\text{ss}}(x) . \end{aligned}$$

The *nonsteady-state addition to free energy* is:

$$\beta^{-1} \gamma(s|\boldsymbol{\mu}, x) \equiv \beta^{-1} \ln \frac{\Pr(\mathcal{S}_t = s | \mathcal{R}_{t-1} \sim \boldsymbol{\mu}, X_t = x)}{\pi_x(s)} .$$

Averaging over observable states this becomes the relative entropy:

$$\langle \gamma(s|\boldsymbol{\mu}, x) \rangle = D_{\text{KL}} [\Pr(\mathcal{S}_t | \mathcal{R}_{t-1} \sim \boldsymbol{\mu}, X_t = x) || \boldsymbol{\pi}_x] ,$$

which is always nonnegative.

Using this setup and decomposing:

$$\frac{\Pr(\mathcal{S}_{0:N} = s^0 \mathbf{s} | \mathcal{R}_{-1} \sim \boldsymbol{\mu}_F, X_{0:N} = x^0 \mathbf{x})}{\Pr(\mathcal{S}_{0:N} = s^{N-1} \mathbf{s}_\leftarrow^R | \mathcal{R}_{-1} \sim \boldsymbol{\mu}_R, X_{0:N} = x^N \mathbf{x}^R)}$$

in analogy with Eq. (21), it is straightforward to extend the remaining results of the main body to the setting in which observed states are functions of a Markov chain. Notably, the path dependencies pick up new contributions from non-Markovity. Also, knowledge of distributions over latent states provides a thermodynamic advantage to Maxwellian Demons.

Appendix B: Integral fluctuation theorems with auxiliary variables

Recall that we quantify how much the auxiliary variable independently informs the state sequence via the nonaveraged conditional mutual information:

$$\begin{aligned} i[\vec{s}; \vec{y} | \vec{x}, \boldsymbol{\mu}_F] &\equiv \ln \frac{\Pr(\vec{s}, \vec{y} | \vec{x}, \boldsymbol{\mu}_F)}{\Pr(\vec{y} | \vec{x}, \boldsymbol{\mu}_F) \Pr(\vec{s} | \vec{x}, \boldsymbol{\mu}_F)} \\ &= \ln \frac{\Pr(\vec{s}, \vec{y}, \vec{x} | \boldsymbol{\mu}_F)}{\Pr(\vec{y}, \vec{x} | \boldsymbol{\mu}_F) \Pr(\vec{s} | \vec{x}, \boldsymbol{\mu}_F)}. \end{aligned}$$

Note that averaging over the input, state, and auxiliary sequences gives the familiar conditional mutual information:

$$\begin{aligned} I[\mathcal{S}_{0:N}; Y_{0:N} | X_{0:N}, \boldsymbol{\mu}_F] \\ = \langle i[\vec{s}; \vec{y} | \vec{x}, \boldsymbol{\mu}_F] \rangle_{\Pr(x_{0:N}, s_{0:N}, y_{0:N} | \boldsymbol{\mu}_F)}. \end{aligned}$$

(Averaging over distributions is the same as being given the distribution, since the distribution over distributions is assumed to be peaked at $\boldsymbol{\mu}_F$.)

Noting that:

$$\begin{aligned} e^{\beta W_{\text{diss}} + i(\vec{s}; \vec{y} | \vec{x}, \boldsymbol{\mu}_F) + \Psi} \\ = e^{\Omega + i(\vec{s}; \vec{y} | \vec{x}, \boldsymbol{\mu}_F) + \Psi + (\gamma_F - \gamma_R)} \\ = \frac{\Pr(\vec{s}, \vec{y}, \vec{x} | \boldsymbol{\mu}_F)}{\Pr(\vec{y}, \vec{x} | \boldsymbol{\mu}_F) \Pr(s^N - 1 \mathbf{s}_L^R | \mathbf{x}^R x^0, \boldsymbol{\mu}_R)} \\ = \frac{\Pr(\vec{s}, \vec{y}, \vec{x} | \boldsymbol{\mu}_F)}{\Pr(\vec{y}, \vec{x} | \boldsymbol{\mu}_F) \Pr(\overleftarrow{s} | \overleftarrow{x}, \boldsymbol{\mu}_R)}, \end{aligned}$$

where $\boldsymbol{\mu}_R = \boldsymbol{\mu}(\boldsymbol{\mu}_F, \vec{x})$, we have the integral fluctuation theorem (IFT):

$$\begin{aligned} &\left\langle e^{-\beta W_{\text{diss}} - i(\vec{s}; \vec{y} | \vec{x}, \boldsymbol{\mu}_F) - \Psi} \right\rangle_{\Pr(\vec{s}, \vec{y}, \vec{x} | \boldsymbol{\mu}_F)} \\ &= \sum_{\vec{x}, \vec{s}, \vec{y}} \Pr(\vec{s}, \vec{y}, \vec{x} | \boldsymbol{\mu}_F) \frac{\Pr(\vec{y}, \vec{x} | \boldsymbol{\mu}_F) \Pr(\overleftarrow{s} | \overleftarrow{x}, \boldsymbol{\mu}_R)}{\Pr(\vec{s}, \vec{y}, \vec{x} | \boldsymbol{\mu}_F)} \\ &= \sum_{\vec{x}, \vec{s}, \vec{y}} \Pr(\vec{y}, \vec{x} | \boldsymbol{\mu}_F) \Pr(\overleftarrow{s} | \overleftarrow{x}, \boldsymbol{\mu}_R) \\ &= \sum_{\vec{x}, \vec{y}} \Pr(\vec{y}, \vec{x} | \boldsymbol{\mu}_F) \sum_{\overleftarrow{s}} \Pr(\overleftarrow{s} | \overleftarrow{x}, \boldsymbol{\mu}_R) \\ &= \sum_{\vec{x}, \vec{y}} \Pr(\vec{y}, \vec{x} | \boldsymbol{\mu}_F) \\ &= 1. \end{aligned}$$

Notably, this relation holds arbitrarily far from equilibrium and allows for the starting and ending distributions to both be nonsteady-state. It is tempting to conclude that the revised Second Law of Thermodynamics should read:

$$\langle W_{\text{diss}} \rangle \geq -k_B T I[\vec{\mathcal{S}}; \vec{Y} | \vec{X}, \boldsymbol{\mu}_F] - \langle Q_{\text{hk}} \rangle, \quad (\text{B1})$$

which includes the effects of both irreversibility and conditional mutual information between state-sequence and auxiliary sequence, given input-sequence. However, we expect that $\langle Q_{\text{hk}} \rangle > 0$, so Eq. (B1) is not the strongest bound derivable. Dropping Ψ from the IFT still yields a true equality, but the derivation runs differently since it depends on the normalization of the conjugate dynamic. Although IFTs with Ψ may be useful for other reasons, it is the non- Ψ IFTs that seems to yield the tighter bound for the revised Second Laws of information thermodynamics without detailed balance.

[1] G. E. Crooks. On thermodynamic and microscopic reversibility. *J. Stat. Mech.: Th. Exp.*, 2011(07):P07008, 2011.

[2] G. E. Crooks. Nonequilibrium measurements of free energy differences for microscopically reversible Markovian systems. *J. Stat. Phys.*, 90(5/6):1481–1487, 1998.

- [3] T. Sagawa and M. Ueda. Nonequilibrium thermodynamics of feedback control. *Phys. Rev. E*, 85:021104, Feb 2012.
- [4] H. Wang and G. Oster. Energy transduction in the F1 motor of ATP synthase. *Nature*, 396(6708):279–282, 1998.
- [5] M. Polettini and M. Esposito. Irreversible thermodynamics of open chemical networks. I. Emergent cycles and broken conservation laws. *J. Chem. Physics*, 141(2), 2014.
- [6] R. Landauer. Statistical physics of machinery: Forgotten middle-ground. *Physica A: Stat. Mech. App.*, 194(1-4):551–562, 1993.
- [7] H. Qian. Nonequilibrium steady-state circulation and heat dissipation functional. *Phys. Rev. E*, 64:022101, 2001.
- [8] W. Horsthemke. Noise induced transitions. In C. Vidal and A. Pacault, editors, *Non-Equilibrium Dynamics in Chemical Systems: Proceedings of the International Symposium, Bordeaux, France, September 3–7, 1984*, pages 150–160, Berlin, Heidelberg, 1984. Springer.
- [9] B. Lindner, J. Garcia-Ojalvo, A. Neiman, and L. Schimansky-Geier. Effects of noise in excitable systems. *Physics Reports*, 392(6):321 – 424, 2004.
- [10] J. P. Crutchfield and C. Aghamohammadi. Not all fluctuations are created equal: Spontaneous variations in thermodynamic function. 2016. Santa Fe Institute Working Paper 16-09-018; arxiv.org:1609.02519 [math-ph].
- [11] U. Seifert. Stochastic thermodynamics, fluctuation theorems and molecular machines. *Rep. Prog. Physics*, 75(12):126001, 2012.
- [12] R. Spinney and I. Ford. Fluctuation relations: A pedagogical overview. In *Nonequilibrium Statistical Physics of Small Systems*, pages 3–56. Wiley-VCH Verlag GmbH & Co. KGaA, 2013.
- [13] Y. Oono and M. Paniconi. Steady state thermodynamics. *Prog. Theo. Phys. Supp.*, 130:29–44, 1998.
- [14] T. Hatano and S. Sasa. Steady-state thermodynamics of Langevin systems. *Phys. Rev. Lett.*, 86:3463–3466, 2001.
- [15] E. H. Trepagnier, C. Jarzynski, F. Ritort, G. E. Crooks, C. J. Bustamante, and J. Liphardt. Experimental test of Hatano and Sasa’s nonequilibrium steady-state equality. *Proc. Natl. Acad. Sci. USA*, 101(42):15038–15041, 2004.
- [16] D. Mandal and C. Jarzynski. Analysis of slow transitions between nonequilibrium steady states. *J. Stat. Mech.: Th. Exp.*, 2016(6):063204, 2016.
- [17] G. E. Crooks. Entropy production fluctuation theorem and the nonequilibrium work relation for free energy differences. *Phys. Rev. E*, 60(3):2721–2726, 1999.
- [18] J. L. England. Statistical physics of self-replication. *J. Chem. Physics*, 139(12):–, 2013.
- [19] J. P. Crutchfield. The calculi of emergence: Computation, dynamics, and induction. *Physica D*, 75:11–54, 1994.
- [20] J. M. Horowitz and S. Vaikuntanathan. Nonequilibrium detailed fluctuation theorem for repeated discrete feedback. *Phys. Rev. E*, 82:061120, Dec 2010.
- [21] We ignore nonergodicity to simplify the development. The approach, though, handles nonergodicity just as well. However, distracting nuances arise that we do not wish to dwell on. For example, if the Markov chain has more than one attracting component for a particular x , then π_x is not unique, but can be constructed as any one of infinitely many probability-normalized linear superpositions of left eigenvectors of $\mathsf{T}^{(\mathcal{S} \rightarrow \mathcal{S}|x)}$ associated with the eigenvalue of unity.
- [22] We start in a discrete-time setup, but later translate to continuous time.
- [23] M. Esposito and C. Van den Broeck. Three detailed fluctuation theorems. *Phys. Rev. Lett.*, 104:090601, Mar 2010.
- [24] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. Wiley-Interscience, New York, second edition, 2006.
- [25] S. Still, D. A. Sivak, A. J. Bell, and G. E. Crooks. Thermodynamics of prediction. *Phys. Rev. Lett.*, 109:120604, Sep 2012.
- [26] J. P. Crutchfield and K. Young. Inferring statistical complexity. *Phys. Rev. Lett.*, 63:105–108, 1989.
- [27] M. Esposito, U. Harbola, and S. Mukamel. Entropy fluctuation theorems in driven open systems: Application to electron counting statistics. *Phys. Rev. E*, 76:031132, Sep 2007.
- [28] G. B. Bagci, U. Tirnakli, and J. Kurths. The second law for the transitions between the non-equilibrium steady states. *Phys. Rev. E*, 87:032161, 2013.
- [29] To be more precise, we write $\Pr(\mathcal{S}_t = s | \mathcal{S}_0 \sim \mu_0, x_{1:t+1})$ as $\Pr_{\mathcal{S}_0 \sim \mu_0}(\mathcal{S}_t = s | x_{1:t+1})$, since the probability is not conditioned on μ_0 —a probability measure for subsequent state sequences. Here, we simply gloss over this nuance, later adopting the shorthand: $\Pr(\mathcal{S}_t = s | \mu_0, x_{1:t+1})$.
- [30] B. Gaveau and L. S. Schulman. A general framework for non-equilibrium phenomena: The master equation and its formal consequences. *Phys. Lett. A*, 229(6):347–353, 1997.
- [31] D. A. Sivak and G. E. Crooks. Near-equilibrium measurements of nonequilibrium free energy. *Phys. Rev. Lett.*, 108(15), 2012.
- [32] H. Qian. Cycle kinetics, steady state thermodynamics and motors: A paradigm for living matter physics. *J. Physics: Cond. Matt.*, 17(47):S3783, 2005.
- [33] Liepelt, S. and Lipowsky, R. Steady-state balance conditions for molecular motor cycles and stochastic nonequilibrium processes. *Euro. Phys. Lett.*, 77(5):50002, 2007.
- [34] S. Liepelt and R. Lipowsky. Kinesin’s network of chemo-mechanical motor cycles. *Phys. Rev. Lett.*, 98:258102, Jun 2007.
- [35] G. E. Crooks. Path-ensemble averages in systems driven far from equilibrium. *Phys. Rev. E*, 61(3):2361–2366, 2000.
- [36] V. Y. Chernyak, M. Chertkov, and C. Jarzynski. Path-integral analysis of fluctuation theorems for general Langevin processes. *J. Stat. Mech.: Th. Exp.*, 2006(08):P08001, 2006.
- [37] R. J. Harris and G. M. Schutz. Fluctuation theorems for stochastic dynamics. *J. Stat. Mech.: Th. Exp.*,

- 2007(07):P07020, 2007.
- [38] U. Seifert. Entropy production along a stochastic trajectory and an integral fluctuation theorem. *Phys. Rev. Lett.*, 95:040602, Jul 2005.
- [39] S. Lahiri and A. M. Jayannavar. Fluctuation theorems for excess and housekeeping heat for underdamped Langevin systems. *Euro. Phys. J. B*, 87(9), 2014.
- [40] C. Jarzynski. Nonequilibrium equality for free energy differences. *Phys. Rev. Lett.*, 78(14):2690–2693, 1997.
- [41] T. Speck and U. Seifert. Integral fluctuation theorem for the housekeeping heat. *J. Phys. A: Math. Gen.*, 38(34):L581, 2005.
- [42] S. Vaikuntanathan and C. Jarzynski. Dissipation and lag in irreversible processes. *EPL (Europhysics Letters)*, 87(6):60005, 2009.
- [43] P. Sartori, L. Granger, C. F. Lee, and J. M. Horowitz. Thermodynamic costs of information processing in sensory adaptation. *PLoS Comput Biol*, 10(12):e1003974, 12 2014.
- [44] E. M. Izhikevich. *Dynamical Systems in Neuroscience*. Computational Neuroscience Series. MIT Press, Boston, Massachusetts, 2010.
- [45] F. Rieke, D. Warland, R. de Ruyter van Steveninck, and W. Bialek. *Spikes: Exploring the Neural Code*. Bradford Books, New York, 1999.
- [46] A. L. Hodgkin and A. F. Huxley. A quantitative description of membrane current and its application to conduction and excitation in nerve. *J. Physio.*, 117(4):500, 1952.
- [47] P. Dayan and L. F. Abbott. *Theoretical Neuroscience: Computational and Mathematical Modeling of Neural Systems*. Computational Neuroscience Series. MIT Press, Boston, Massachusetts, revised edition, 2005.
- [48] J. Patlak. Molecular kinetics of voltage-dependent Na^+ channels. *Physiol. Rev.*, 71(4):1047–1080, 1991.
- [49] J. P. Crutchfield, C. J. Ellison, and P. M. Riechers. Exact complexity: Spectral decomposition of intrinsic computation. *Phys. Lett. A*, 380(9-10):998–1002, 2016.
- [50] P. M. Riechers and J. P. Crutchfield. Beyond the spectral theorem: Decomposing arbitrary functions of nondiagonalizable operators. 2016. arxiv.org:1607.06526 [math-ph].
- [51] The characteristic timescale is actually the net result of a combination of timescales from the inverse eigenvalues of G . Of necessity, these are the same timescales that determine the relaxation of the state distribution.
- [52] D. Lacoste, A. W. C. Lau, and K. Mallick. Fluctuation theorem and large deviation function for a solvable model of a molecular motor. *Phys. Rev. E*, 78:011915, 2008.
- [53] Y. Murashita, K. Funo, and M. Ueda. Nonequilibrium equalities in absolutely irreversible processes. *Phys. Rev. E*, 90:042110, Oct 2014.
- [54] B. Altaner, A. Wachtel, and J. Vollmer. Fluctuating currents in stochastic thermodynamics II: Energy conversion and nonequilibrium response in kinesin models. *arXiv:1504.03648 [cond-mat.stat-mech]*, 2015.
- [55] D. Colquhoun and A. G. Hawkes. Relaxation and fluctuations of membrane currents that flow through drug-operated channels. *Proc. Roy. Soc. Lond. B: Bio. Sci.*, 199(1135):231–262, 1977.
- [56] S. Lahiri and S. Ganguli. A memory frontier for complex synapses. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Adv. Neural Info. Proc. Sys. 26*, pages 1034–1042. Curran Associates, Inc., 2013.
- [57] Q. Shen, Q. Hao, and S. M. Gruner. Macromolecular phasing. *Physics Today*, 59(3):46–52, 2006.
- [58] P. W. Anderson. More is different. *Science*, 177(4047):393–396, 1972.
- [59] A. B. Boyd, D. Mandal, and J. P. Crutchfield. Correlation-powered information engines and the thermodynamics of self-correction. 2016. arXiv.org:1606.08506 [cond-mat.stat-mech].
- [60] A. B. Boyd, D. Mandal, and J. P. Crutchfield. Leveraging environmental correlations: The thermodynamics of requisite variety. 2016. arXiv.org:1609.05353 [cond-mat.stat-mech].