# Exact fluctuation relation for open systems beyond the Jarzynski equality

Mohammad Rahbar

*Technical University of Munich; TUM School of Natural Sciences,*
*Department of Chemistry, Lichtenbergstr. 4, D-85748 Garching, Germany*

Christopher J. Stein

*Technical University of Munich; TUM School of Natural Sciences,*
*Department of Chemistry, Catalysis Research Center,*
*Atomistic Modeling Center, Munich Data Science Institute,*
*Lichtenbergstr. 4, D-85748 Garching, Germany**

(Dated: November 14, 2025)

We derive exact fluctuation equalities for open systems that recover free energy differences between two equilibrium endpoints connected by nonequilibrium processes with arbitrary dynamics and coupling. The exponential of the free energy difference is expressed in terms of ensemble averages of the Hamiltonian of mean force (HMF) shift and the chi-squared divergence between the initial and final marginal probability distribution of the open system. A trajectory counterpart of this relation follows from an asymptotic equilibration postulate, which treats relaxation to the final stationary canonical state as a boundary condition rather than as a consequence of constraints on the driven dynamics. In the frozen-coupling regime, the HMF shift reduces to the bare-system Hamiltonian shift, yielding a clear heat–work decomposition. The Jarzynski equality (JE) is recovered under the assumption of Hamiltonian dynamics for the combined system. We validate the theory on a dissipative, phase-space-compressing drive followed by an underdamped Langevin relaxation, where the assumptions underlying the JE break down, whereas our equality reproduces the exact free energy differences.

The Jarzynski equality (JE) [1], $\langle e^{-\beta W} \rangle = e^{-\beta \Delta F}$, is a fundamental result in modern statistical mechanics. It establishes an exact connection between the exponential average of the nonequilibrium work $W$ performed during an irreversible process and the equilibrium free energy difference $\Delta F$. Along with related fluctuation theorems (FTs) [2–4], it has been confirmed in both classical and quantum systems and provides a powerful framework for extracting equilibrium information from nonequilibrium measurements [5–12]. Nevertheless, its generality must be interpreted with care in different situations [13–28]. The theoretical foundation of the original JE rests on strict dynamical reversibility [1]. In its original derivation, the composite system $\mathcal{S} + \mathcal{E}$, with bare system $\mathcal{S}$ and environment $\mathcal{E}$, evolves under deterministic dynamics that obey Liouville's theorem, ensuring phase–space volume preservation. This preservation implies a fundamental dynamical constraint underlying the equality. A stochastic generalization soon followed, reformulating the dynamics through a Markov master equation or Langevin description [29]. In this generalization, reversibility takes a statistical form through detailed balance (DB). DB ensures that, for each fixed control parameter $\lambda$, the Gibbs–Boltzmann canonical distribution serves as the stationary state probability distribution, providing the probabilistic analogue of microscopic reversibility. The necessity of reversibility constraints was highlighted by Cohen and Mauzerall's criticism [22], which questioned the formal validity of the JE

and suggested its applicability might be limited to near-equilibrium weakly coupled regimes. Jarzynski's reply [30] demonstrated that the JE is a mathematical identity following directly from Hamiltonian dynamics of the composite system ($\mathcal{S} + \mathcal{E}$). This derivation, later refined to incorporate strong coupling through the Hamiltonian of mean force (HMF) $H_\beta^*(X_\mathcal{S}, \lambda)$, confirmed that while the free energy interpretation changes, the proof still relies on Liouville's theorem and therefore on microscopic reversibility of the composite dynamics. Speck and Seifert [31] extended the equality to non-Markovian processes governed by the generalized Langevin equation. For such dynamics, the JE remains valid, provided the friction kernel and noise correlations obey the fluctuation–dissipation theorem (FDT). Compliance with FDT ensures that the Gibbs–Boltzmann distribution remains stationary for each frozen $\lambda$. While the JE demonstrated resilience against classical memory effects, a more fundamental challenge arose from an observation that the standard JE is necessarily violated in the presence of feedback — a phenomenon central to Maxwell's Demon [32] and experimentally realized in information ratchets[33]. A formal resolution was achieved through the generalized Jarzynski equality (GJE) under feedback control, as formulated by Sagawa and Ueda [14, 21]. That extension incorporates the stochastic mutual information gained during measurement as a thermodynamic resource. However, this GJE still relies on a dynamical assumption, requiring the classical stochastic processes to satisfy local detailed balance (LDB) to ensure that the time-reversal property of path probabilities is well defined [34]. Regardless of the distinct approaches for the JE derivations

* christopher.stein@tum.de

and related generalizations, they all share an inherent dependency on underlying dynamical constraints such as microscopic reversibility, DB, FDT, or LDB [2, 35]. When these constraints are violated, the JE (or GJE) is no longer guaranteed to hold. Across biology, chemistry, physics, and engineering, many systems are driven under conditions that break these dynamical constraints, yet still relax to a genuine equilibrium once the driving is kept fixed. [36–48].

The key question is whether it is possible to quantify the free energy difference between two equilibrium points connected by nonequilibrium processes, given arbitrary dynamics and system-environment coupling. In this Letter, we derive equalities for this most general scenario. As the first result, we derive an endpoint equality in terms of the driving and coupling protocols (Thm. 1), showing that its validity is independent of both the underlying dynamics and the strength of system–environment coupling linking the equilibrium endpoints. The second result (Thm. 2) establishes the trajectory counterpart based on the asymptotic postulate of thermal equilibration. For strongly coupled systems in the frozen–coupling regime — the physically most natural case — the equalities depend on the HMF in exactly the same way as on the bare–system Hamiltonian (Cor. 1). We then derive a heat–work decomposition (Cor. 2), which reveals the thermodynamic structure of the relations and provides a consistent interpretation of work and heat exchange. In the limit of Liouvillian dynamics and frozen coupling, our framework reduces to the standard JE. We validate the theory on a composite model, where the standard work estimator of JE fails under non-Liouvillian driving, while the endpoint and trajectory counterpart equalities reproduce the exact equilibrium free energy differences. We begin with a few necessary definitions (Def. 1-3) that later serve as the foundation for the endpoint and trajectory relations (Thm. 1 & 2) derived in this work.

*Definition* 1 (Microscopic trajectory of the composite system). Here, we formulate the evolution of the composite system in an abstract manner, remaining agnostic to the details of its underlying dynamics. We consider a composite system $\mathcal{S} + \mathcal{E}$, consisting of a system of interest $\mathcal{S}$ and its environment $\mathcal{E}$. The total phase space of the composite at time $t$ is denoted by $\Gamma_t$, and a single microscopic state by $X \in \Gamma_t$. Each microstate contains the full phase-space coordinates of both parts, $X = (X_{\mathcal{S}}, X_{\mathcal{E}})$, where $X_{\mathcal{S}}$ and $X_{\mathcal{E}}$ collect the positions and momenta of particles in $\mathcal{S}$ and $\mathcal{E}$, respectively. We introduce the trajectory map $\mathcal{T}_t : \Gamma_0 \to \Gamma_t$, which assigns to each initial state $X_0 \in \Gamma_0$, its evolved image $X(t|X_0) = \mathcal{T}_t(X_0)$ at time $t$. The map $\mathcal{T}_t$ is a purely kinematic construct and does not rely on any specific dynamical generator. It may represent Hamiltonian, stochastic, or arbitrary evolution, providing a unified description for all possible dynamics. Since $\mathcal{T}_t$ acts on the composite microstates $X_0 \in \Gamma_0$, no assumption is made that the system and environment trajectories can be defined independently — a requirement intrinsic to any open
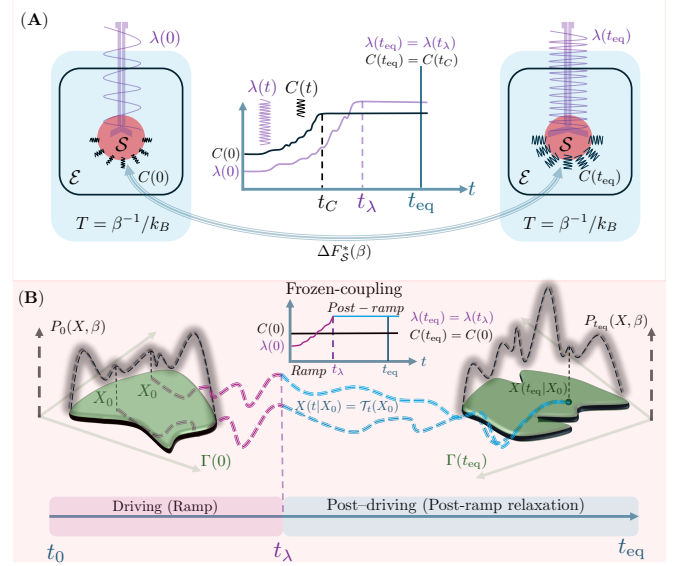


FIG. 1. **(A)** Schematic of the composite $\mathcal{S} + \mathcal{E}$ at a fixed bath temperature $T = \beta^{-1}/k_B$. The control $\lambda(t)$ (schematic purple spring) and the coupling $C(t)$ (schematic black spring) are driven independently and become constant at times $t_\lambda$ and $t_C$, respectively; the composite system then relaxes to equilibrium at $t_{\rm eq}$. The system free energy difference $\Delta F_{\mathcal{S}}^*(\beta)$ refers to the canonical endpoints $(\lambda(0), C(0)) \to (\lambda(t_{\rm eq}), C(t_{\rm eq}))$ and is defined through the HMF partition functions. **(B)** An initial microstate $X_0 \sim P_0(X, \beta)$ on $\Gamma(0)$ evolves under the kinematic map $X(t|X_0) = \mathcal{T}_t(X_0)$ through the ramp and the post–ramp relaxation, yielding the asymptotic density $P_{t_{\rm eq}}(X, \beta)$ on $\Gamma(t_{\rm eq})$. The frozen coupling case, $C(t) \equiv C(0)$, is the regime used in our validation.

system description. The microscopic state of the composite system along a trajectory is therefore written as $X(t|X_0) = \big(X_{\mathcal{S}}(t|X_0) = \mathcal{T}_t^{\mathcal{S}}(X_0), X_{\mathcal{E}}(t|X_0) = \mathcal{T}_t^{\mathcal{E}}(X_0)\big)$, where $\mathcal{T}_t^{\mathcal{S}}$ and $\mathcal{T}_t^{\mathcal{E}}$ denote the respective projections of the composite map $\mathcal{T}_t$ onto the system and environment coordinates. Next, we specify the Hamiltonian structure along the established trajectory.

*Definition* 2 (Total Hamiltonian along a trajectory). The instantaneous total Hamiltonian of the composite system $\mathcal{S} + \mathcal{E}$ evaluated along the trajectory $\mathcal{T}_t(X_0)$ is defined as

$$\mathcal{H}_{\mathcal{S}+\mathcal{E}}\big(\mathcal{T}_t(X_0), \lambda(t), C(t)\big) = \mathcal{H}_{\mathcal{S}}\big(\mathcal{T}_t^{\mathcal{S}}(X_0), \lambda(t)\big) + \mathcal{H}_{\mathcal{E}}\big(\mathcal{T}_t^{\mathcal{E}}(X_0)\big) + \mathcal{V}_{\mathcal{S}\mathcal{E}}\big(\mathcal{T}_t^{\mathcal{S}}(X_0), \mathcal{T}_t^{\mathcal{E}}(X_0), C(t)\big). \quad (1)$$

Here $\mathcal{H}_{\mathcal{S}}$ and $\mathcal{H}_{\mathcal{E}}$ are the bare Hamiltonians of the system and environment, respectively, and $\mathcal{V}_{\mathcal{S}\mathcal{E}}$ is their interaction potential. Two externally controlled parameters appear: $\lambda(t)$ represents the driving protocol acting on $\mathcal{S}$ (for example, a control parameter or mechanical coordinate) and $C(t)$ controls the coupling strength or interaction channel between $\mathcal{S}$ and $\mathcal{E}$. By allowing $\lambda(t)$ and $C(t)$ to vary independently, we retain the most general energetic structure of an open system subject to simultaneously applied driving and coupling controls. Each

protocol may correspond to a sudden quench or to a continuous ramp in time. We denote by $t_\lambda$ and $t_c$ the times at which the corresponding protocols become constant, $\lambda(t) = \text{const}$ for $t \geq t_\lambda$ and $C(t) = \text{const}$ for $t \geq t_c$. The composite system subsequently relaxes to equilibrium at a later time $t_{\text{eq}}$. In general, $0 \leq t_\lambda,\ t_c \leq t_{\text{eq}}$ such that $t_\lambda \neq t_c$, so the completion of driving, coupling, and equilibration do not necessarily coincide. This distinction will later be crucial when we analyze the frozen–coupling limit $C(t) \equiv C(0)$. Now, we introduce the canonical ensembles for the composite system and the system of interest $\mathcal{S}$. These equilibrium measures serve as statistical reference points for the initial and final endpoints of the process.

*Definition* 3 (Canonical ensembles of the composite system and marginal distributions of the system). The equilibrium probability distribution of the composite system $\mathcal{S} + \mathcal{E}$ at temperature $\beta^{-1}/k_{\text{B}}$ is defined in canonical form as $P(X, \beta) = e^{-\beta\,\mathcal{H}_{\mathcal{S}+\mathcal{E}}(X, \lambda, C)}/\mathcal{Z}_{\mathcal{S}+\mathcal{E}}(\lambda, C, \beta)$, where $\mathcal{H}_{\mathcal{S}+\mathcal{E}}$ is the total Hamiltonian given in Eq. (1) and $\mathcal{Z}_{\mathcal{S}+\mathcal{E}}$ denotes the partition function of the composite system, $\mathcal{Z}_{\mathcal{S}+\mathcal{E}}(\lambda, C, \beta) = \int_{\Gamma_t} dX\, e^{-\beta\,\mathcal{H}_{\mathcal{S}+\mathcal{E}}(X, \lambda, C)}$. Two specific equilibrium ensembles are relevant to the present discussion, i.e.

$$P_0(X, \beta) = \frac{e^{-\beta\,\mathcal{H}_{\mathcal{S}+\mathcal{E}}(X, \lambda(0), C(0))}}{\mathcal{Z}_{\mathcal{S}+\mathcal{E}}(\lambda(0), C(0), \beta)}, \tag{2}$$

$$P_{t_{\text{eq}}}(X, \beta) = \frac{e^{-\beta\,\mathcal{H}_{\mathcal{S}+\mathcal{E}}(X, \lambda(t_{\text{eq}}), C(t_{\text{eq}}))}}{\mathcal{Z}_{\mathcal{S}+\mathcal{E}}(\lambda(t_{\text{eq}}), C(t_{\text{eq}}), \beta)}. \tag{3}$$

The distributions $P_0$ and $P_{t_{\text{eq}}}$ describe the initial and final equilibrium states of the composite system at the same bath temperature $\beta^{-1}/k_{\text{B}}$. The equilibrium marginal distribution of the system $\mathcal{S}$ is obtained by integrating over the environmental coordinates $P^{\mathcal{S}}(X_{\mathcal{S}}, \beta) = \int dX_{\mathcal{E}}\, P(X_{\mathcal{S}}, X_{\mathcal{E}}, \beta)$. This probability can be expressed in canonical form through the HMF, defined as

$$\mathcal{H}_\beta^*(X_{\mathcal{S}}, \lambda, C) = \mathcal{H}_{\mathcal{S}}(X_{\mathcal{S}}, \lambda)$$
$$- \frac{1}{\beta} \ln \int dX_{\mathcal{E}}\, e^{-\beta\,[\mathcal{H}_{\mathcal{E}}(X_{\mathcal{E}}) + \mathcal{V}_{\mathcal{S}\mathcal{E}}(X_{\mathcal{S}}, X_{\mathcal{E}}, C)]}. \tag{4}$$

Equation (4) defines an effective Hamiltonian that incorporates the influence of the environment on $\mathcal{S}$ at fixed $(\lambda, C, \beta)$. The corresponding partition function of the system is $\mathcal{Z}_{\mathcal{S}}^*(\lambda, C, \beta) = \int dX_{\mathcal{S}}\, e^{-\beta\,\mathcal{H}_\beta^*(X_{\mathcal{S}}, \lambda, C)}$. Hence the equilibrium marginal distribution of the system reads $P^{\mathcal{S}}(X_{\mathcal{S}}, \beta) = e^{-\beta\,\mathcal{H}_\beta^*(X_{\mathcal{S}}, \lambda, C)}/\mathcal{Z}_{\mathcal{S}}^*(\lambda, C, \beta)$. Applying them to the two equilibrium endpoints of the process yields

$$P_0^{\mathcal{S}}(X_{\mathcal{S}}, \beta) = \frac{e^{-\beta\,\mathcal{H}_\beta^*(X_{\mathcal{S}}, \lambda(0), C(0))}}{\mathcal{Z}_{\mathcal{S}}^*(\lambda(0), C(0), \beta)}, \tag{5}$$

$$P_{t_{\text{eq}}}^{\mathcal{S}}(X_{\mathcal{S}}, \beta) = \frac{e^{-\beta\,\mathcal{H}_\beta^*(X_{\mathcal{S}}, \lambda(t_{\text{eq}}), C(t_{\text{eq}}))}}{\mathcal{Z}_{\mathcal{S}}^*(\lambda(t_{\text{eq}}), C(t_{\text{eq}}), \beta)}. \tag{6}$$

Eqs. (5) and (6) define the initial and final equilibrium probability distributions of the system in terms of its HMF, corresponding to the protocol endpoints $(\lambda(0), C(0))$ and $(\lambda(t_{\text{eq}}), C(t_{\text{eq}}))$, respectively. In Thm. 1, we show the exact connection between the system free energy difference at two endpoints and the canonical distributions and related HMFs at those endpoints.

*Theorem* 1 (Endpoint equalities for free energy differences). For two equilibrium endpoints of the composite system $\mathcal{S} + \mathcal{E}$ prepared at the same inverse temperature $\beta$, we have (see Appendix for the proof)

$$e^{-\beta\Delta F_{\mathcal{S}}^*(\beta)} = \frac{\left\langle e^{-\beta\Delta\mathcal{H}_\beta^*(X_{\mathcal{S}})} \right\rangle_{\mathcal{S}}}{1 + \chi^2\!\left(P_{t_{\text{eq}}}^{\mathcal{S}} \parallel P_0^{\mathcal{S}}\right)}, \tag{7}$$

and

$$e^{+\beta\Delta F_{\mathcal{S}}^*(\beta)} = \left\langle e^{+\beta\Delta\mathcal{H}_\beta^*(X_{\mathcal{S}})} \right\rangle_{\mathcal{S}}. \tag{8}$$

The free energy difference is

$$\Delta F_{\mathcal{S}}^*(\beta) = F_{\mathcal{S}}^*(\lambda(t_{\text{eq}}), C(t_{\text{eq}}), \beta) - F_{\mathcal{S}}^*(\lambda(0), C(0), \beta), \tag{9}$$

where $F_{\mathcal{S}}^*(\lambda, C, \beta) = \beta^{-1} \ln \mathcal{Z}_{\mathcal{S}}^*(\lambda, C, \beta)$ and the HMF shift is

$$\Delta\mathcal{H}_\beta^*(X_{\mathcal{S}}) = \mathcal{H}_\beta^*(X_{\mathcal{S}}, \lambda(t_{\text{eq}}), C(t_{\text{eq}}))$$
$$- \mathcal{H}_\beta^*(X_{\mathcal{S}}, \lambda(0), C(0)). \tag{10}$$

Averages are taken over the final equilibrium ensemble,

$$\langle \bullet \rangle_{\mathcal{S}} = \int dX_{\mathcal{S}}\, \bullet\, P_{t_{\text{eq}}}^{\mathcal{S}}(X_{\mathcal{S}}, \beta), \tag{11}$$

and the chi–squared divergence between the endpoint marginals is

$$1 + \chi^2\!\left(P_{t_{\text{eq}}}^{\mathcal{S}} \parallel P_0^{\mathcal{S}}\right) = \int dX_{\mathcal{S}}\, \frac{\left(P_{t_{\text{eq}}}^{\mathcal{S}}(X_{\mathcal{S}}, \beta)\right)^2}{P_0^{\mathcal{S}}(X_{\mathcal{S}}, \beta)}. \tag{12}$$

The result established here coincides with the endpoint equality derived in [49], where the protocol is implicitly present. As can be seen from the derivation of Eqs. (7) and (8), no restriction is imposed on either the coupling or the form of the dynamics. The reasoning relies only on the existence of well-defined equilibrium endpoints, which guarantees that the corresponding marginal distributions and free energy difference are statistically and thermodynamically meaningful, respectively. Before introducing the trajectory counterparts, it is essential to clarify the assumption underlying their derivation. We postulate *asymptotic equilibration* of the full probability distribution along the trajectory. This assumption expresses the natural requirement that the Helmholtz free energy difference $\Delta F$ be well defined from the trajectory perspective of any nonequilibrium process connecting two equilibrium states, without imposing any particular dynamical constraint during the driving stage. It contrasts sharply with the standard derivations of the JE and its

extensions, where the equality holds only if the underlying dynamics, during the application of the driving protocol, evolve toward a canonical stationary state once the control parameters are held fixed. In other words, the stationary canonical distribution at equilibrium is a consequence of constraints on the driven dynamics. By treating asymptotic convergence to $P_{t_{eq}}$ as a natural boundary condition of thermal contact — rather than as a consequence of these dynamical constraints — the definition of equilibrium becomes decoupled from the specific generator of dynamics during the driving stage. This decoupling extends the application of our derivations to a broader class of systems, including those governed by arbitrary dynamics that are incompatible with the underlying assumptions of JE. We will later discuss the subtleties surrounding the imposition of that boundary condition in the validation section.

*Theorem* 2 (Trajectory counterpart of the endpoint equalities). Using the trajectory map $\mathcal{T}_t : \Gamma_0 \rightarrow \Gamma_t$ defined in Def. 1, the time evolution of the probability density of the composite system can be defined as

$$P_t(X, \beta) = \int_{\Gamma_0} \mathrm{d}X_0\, \delta\big(X - \mathcal{T}_t(X_0)\big)\, P_0(X_0, \beta), \quad (13)$$

where $P_0(X_0, \beta)$ is the initial canonical distribution given in Eq. (2). Integrating over the environmental degrees of freedom yields the time evolution of the system's marginal distribution (see Appendix for derivation)

$$P_t^{\mathcal{S}}(X_{\mathcal{S}}, \beta) = \int_{\Gamma_0} \mathrm{d}X_0\, P_0(X_0, \beta)\, \delta\big(X_{\mathcal{S}} - \mathcal{T}_t^{\mathcal{S}}(X_0)\big). \quad (14)$$

We postulate the asymptotic equilibration of the full distribution that implies convergence of the time evolution of the system's marginal distribution to its canonical endpoint form

$$\lim_{t \to t_{eq}} P_t^{\mathcal{S}}(X_{\mathcal{S}}, \beta) = P_{t_{eq}}^{\mathcal{S}}(X_{\mathcal{S}}, \beta). \quad (15)$$

Under that condition, the ensemble averages in Thm. 1 admit the exact trajectory representation (see Appendix for proof), with averages over the initial canonical ensemble, $\langle \bullet \rangle_{X_0} = \int \mathrm{d}X_0 \, \bullet \, P_0(X_0, \beta)$. Consequently, the endpoint equalities (7) and (8) take the trajectory form

$$e^{-\beta \Delta F_{\mathcal{S}}^*(\beta)} = \frac{\left\langle e^{-\beta \Delta \mathcal{H}_\beta^*\left(\mathcal{T}_{t_{eq}}^{\mathcal{S}}(X_0)\right)} \right\rangle_{X_0}}{1 + \chi^2\big(P_{t_{eq}}^{\mathcal{S}} \parallel P_0^{\mathcal{S}}\big)}, \quad (16)$$

and

$$e^{+\beta \Delta F_{\mathcal{S}}^*(\beta)} = \left\langle e^{+\beta \Delta \mathcal{H}_\beta^*\left(\mathcal{T}_{t_{eq}}^{\mathcal{S}}(X_0)\right)} \right\rangle_{X_0}. \quad (17)$$

*Frozen coupling regime.* The scenario of frozen coupling deserves particular attention. In most free energy difference calculations — both theoretical and experimental — the system remains in fixed coupling (the coupling strength $C$ between $\mathcal{S}$ and $\mathcal{E}$ is kept constant throughout the process) with its environment while the control parameter $\lambda(t)$ is driven. This is precisely the regime underlying the standard derivation of the JE, although the JE itself is agnostic to whether the protocol acts as a driving or a coupling. Within this common, physically motivated freezing scenario, the formal structure of our derivation simplifies considerably.

*Corollary* 1 (Frozen coupling: trajectory forms reduce to bare-system expressions). For a fixed coupling protocol $C(t) \equiv C(0)$, the HMF increment along a trajectory reduces exactly to the difference of bare–system Hamiltonians,

$$\Delta \mathcal{H}_\beta^*\big(\mathcal{T}_{t_{eq}}^{\mathcal{S}}(X_0)\big)$$
$$= \mathcal{H}_\beta^*\big(\mathcal{T}_{t_{eq}}^{\mathcal{S}}(X_0), \lambda(t_{eq}), C(0)\big) - \mathcal{H}_\beta^*\big(\mathcal{T}_{t_{eq}}^{\mathcal{S}}(X_0), \lambda(0), C(0)\big)$$
$$= \mathcal{H}_{\mathcal{S}}\big(\mathcal{T}_{t_{eq}}^{\mathcal{S}}(X_0), \lambda(t_{eq})\big) - \mathcal{H}_{\mathcal{S}}\big(\mathcal{T}_{t_{eq}}^{\mathcal{S}}(X_0), \lambda(0)\big), \quad (18)$$

Equation (18) follows directly from Eq. (4), for constant $C$, the environmental contribution inside the logarithmic term cancels between the two endpoints. Substituting Eq. (18) into the trajectory equalities (16)–(17) yields

$$e^{-\beta \Delta F_{\mathcal{S}}^*(\beta)}$$
$$= \frac{\left\langle e^{-\beta \mathcal{H}_{\mathcal{S}}\left(\mathcal{T}_{t_{eq}}^{\mathcal{S}}(X_0), \lambda(t_{eq})\right)} e^{+\beta \mathcal{H}_{\mathcal{S}}\left(\mathcal{T}_{t_{eq}}^{\mathcal{S}}(X_0), \lambda(0)\right)} \right\rangle_{X_0}}{1 + \chi^2\big(P_{t_{eq}}^{\mathcal{S}} \parallel P_0^{\mathcal{S}}\big)},$$
$$(19)$$

and the complementary identity,

$$e^{+\beta \Delta F_{\mathcal{S}}^*(\beta)}$$
$$= \left\langle e^{+\beta \mathcal{H}_{\mathcal{S}}\left(\mathcal{T}_{t_{eq}}^{\mathcal{S}}(X_0), \lambda(t_{eq})\right)} e^{-\beta \mathcal{H}_{\mathcal{S}}\left(\mathcal{T}_{t_{eq}}^{\mathcal{S}}(X_0), \lambda(0)\right)} \right\rangle_{X_0}.$$
$$(20)$$

It is important to mention that as soon as $C(t)$ is non-frozen, the cancellation in (18) no longer occurs and the full HMF structure must be retained.

*Corollary* 2 (Heat–work decomposition for frozen coupling). We now discuss the thermodynamic structure hidden in Eqs. (19) and (20). In the frozen coupling regimes, energetic changes in $\mathcal{S}$ arise from (i) external manipulation of $\lambda(t)$ and (ii) exchange with the environment through the microscopic evolution of $X_{\mathcal{E}}$. To inspect these contributions, we start from Eq. (18) and apply a simple algebraic insertion–subtraction of the initial energy at fixed protocol, $\mathcal{H}_{\mathcal{S}}\big(X_{\mathcal{S}}(0|X_0), \lambda(0)\big)$, which

yields the following representation. For $C(t) \equiv C(0)$,

$$\mathcal{H}_\beta^*\big(\mathcal{T}_{t_{\mathrm{eq}}}^{\mathcal{S}}(X_0), \lambda(t_{\mathrm{eq}}), C(0)\big) - \mathcal{H}_\beta^*\big(\mathcal{T}_{t_{\mathrm{eq}}}^{\mathcal{S}}(X_0), \lambda(0), C(0)\big)$$

$$= \underbrace{\Big[\mathcal{H}_{\mathcal{S}}\big(X_{\mathcal{S}}(t_{\mathrm{eq}}|X_0), \lambda(t_{\mathrm{eq}})\big) - \mathcal{H}_{\mathcal{S}}\big(X_{\mathcal{S}}(0|X_0), \lambda(0)\big)\Big]}_{I(t_{\mathrm{eq}}|X_0)}$$

$$- \underbrace{\Big[\mathcal{H}_{\mathcal{S}}\big(X_{\mathcal{S}}(t_{\mathrm{eq}}|X_0), \lambda(0)\big) - \mathcal{H}_{\mathcal{S}}\big(X_{\mathcal{S}}(0|X_0), \lambda(0)\big)\Big]}_{II(t_{\mathrm{eq}}|X_0)}.$$

$$(21)$$

*Term I.* Along each realization, the system follows the trajectory $X_{\mathcal{S}}(t|X_0)$ in its phase space. The total time variation of the bare–system Hamiltonian $\mathcal{H}_{\mathcal{S}}(X_{\mathcal{S}}(t|X_0), \lambda(t))$ is obtained by applying the chain rule to its explicit dependence on both $X_{\mathcal{S}}$ and the control parameter $\lambda(t)$,

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathcal{H}_{\mathcal{S}}\big(X_{\mathcal{S}}(t|X_0), \lambda(t)\big) = \underbrace{\big(\nabla_{X_S}\mathcal{H}_{\mathcal{S}}\big) \cdot \dot{X}_{\mathcal{S}}}_{\dot{Q}_{\mathcal{S}}} + \underbrace{\frac{\partial \mathcal{H}_{\mathcal{S}}}{\partial \lambda}\dot{\lambda}}_{\dot{W}_{\mathcal{S}}}.$$

$$(22)$$

The first term describes the energy exchange with the environment along the actual phase–space trajectory and defines the instantaneous heat rate, $\dot{Q}_{\mathcal{S}}$, whereas the second term accounts for the parametric energy input due to the protocol and defines the instantaneous work rate, $\dot{W}_{\mathcal{S}}$. Following the standard conventions of stochastic energetics [50], their time integrals give

$$W_{\mathcal{S}}(t_{\mathrm{eq}}|X_0) = \int_0^{t_{\mathrm{eq}}} \mathrm{d}t\, \frac{\partial \mathcal{H}_{\mathcal{S}}}{\partial \lambda}\dot{\lambda}(t), \qquad (23)$$

$$Q_{\mathcal{S}}(t_{\mathrm{eq}}|X_0) = \int_0^{t_{\mathrm{eq}}} \mathrm{d}t\, \big(\nabla_{X_S}\mathcal{H}_{\mathcal{S}}\big) \cdot \dot{X}_{\mathcal{S}}(t|X_0), \qquad (24)$$

so that integrating Eq. (22) over $t \in [0, t_{\mathrm{eq}}]$ yields the total increment

$$I(t_{\mathrm{eq}}|X_0) = W_{\mathcal{S}}(t_{\mathrm{eq}}|X_0) + Q_{\mathcal{S}}(t_{\mathrm{eq}}|X_0). \qquad (25)$$

*Term II.* For a fixed control $\lambda(0)$, $II(t_{\mathrm{eq}}|X_0)$ can be written as

$$II(t_{\mathrm{eq}}|X_0) = \int_0^{t_{\mathrm{eq}}} \mathrm{d}t\, \big(\nabla_{X_S}\mathcal{H}_{\mathcal{S}}(X_{\mathcal{S}}, \lambda(0))\big) \cdot \dot{X}_{\mathcal{S}}. \quad (26)$$

Here the gradient is evaluated with respect to $\mathcal{H}_{\mathcal{S}}$ frozen at $\lambda(0)$, whereas the $\dot{X}_{\mathcal{S}}$ belongs to the driven trajectory generated by $\lambda(t)$. Consequently, $II$ is not a thermodynamic heat unless $\dot{\lambda} = 0$ for all $t$; it is a reference functional that projects the driven generalized velocity $(\dot{X}_{\mathcal{S}})$ onto the generalized force field $(\nabla_{X_S}\mathcal{H}_{\mathcal{S}})$ of the initial stationary protocol. Because the reference and realized dynamics belong to distinct protocol layers, $II$

acts as a feedback-like correction without a definite thermodynamic notion. Combining Eqs. (25) and (26) with Eq. (21), the total bare–system increment becomes

$$W_{\mathcal{S}}(t_{\mathrm{eq}}|X_0) + Q_{\mathcal{S}}(t_{\mathrm{eq}}|X_0) - II(t_{\mathrm{eq}}|X_0). \qquad (27)$$

Substituting Eq. (27) into Eqs. (19)–(20) yields the compact heat–work representation of the endpoint equalities,

$$e^{-\beta \Delta F_{\mathcal{S}}^*(\beta)} = \frac{\big\langle e^{-\beta[W_{\mathcal{S}} + Q_{\mathcal{S}} - II]}\big\rangle_{X_0}}{1 + \chi^2(P_{t_{\mathrm{eq}}}^{\mathcal{S}} \parallel P_0^{\mathcal{S}})}, \qquad (28)$$

$$e^{+\beta \Delta F_{\mathcal{S}}^*(\beta)} = \big\langle e^{+\beta[W_{\mathcal{S}} + Q_{\mathcal{S}} - II]}\big\rangle_{X_0}. \qquad (29)$$

The decomposition isolates three pathwise contributions: (i) $W_{\mathcal{S}}$, mechanical work due to driving; (ii) $Q_{\mathcal{S}}$, total heat exchanged during the whole process; and (iii) $II$, a reference (projection) functional obtained by contracting the driven velocity $(\dot{X}_{\mathcal{S}})$ with the $\lambda(0)$ force field $(\nabla_{X_S}\mathcal{H}_{\mathcal{S}})$. Only when $\lambda$ is held fixed does $II$ coincide with a thermodynamic heat.

*GJE connection.* The GJE restores $e^{-\beta \Delta F}$ under feedback by augmenting work with mutual information, providing an operational measure of information as a thermodynamic resource [21]. Our identities achieve restoration without invoking measurement–feedback a priori and without assuming LDB. The pathwise correction $Q_{\mathcal{S}} - II$ compensates dynamical asymmetry during driving, and the ensemble overlap $1 + \chi^2(P_{t_{\mathrm{eq}}}^{\mathcal{S}} \parallel P_0^{\mathcal{S}})$ compensates endpoint mismatch. Formally, our equality is broader in scope, it holds without feedback and without LDB, while its structure is fully compatible with the GJE.

*JE as limiting case.* The central identity underlying all our derivations is

$$e^{-\beta \Delta F_{\mathcal{S}}^*} = \frac{Z^*(\lambda(t_{\mathrm{eq}}), C(t_{\mathrm{eq}}), \beta)}{Z^*(\lambda(0), C(0), \beta)}. \qquad (30)$$

By the definition of the HMF, $Z_{\mathcal{S}+\mathcal{E}}(\lambda, C, \beta) = Z^*(\lambda, C, \beta)\, Z_{\mathcal{E}}(\beta)$, so Eq. (30) is equivalent to

$$\frac{Z^*(\lambda(t_{\mathrm{eq}}), C(t_{\mathrm{eq}}), \beta)}{Z^*(\lambda(0), C(0), \beta)} = \frac{Z_{\mathcal{S}+\mathcal{E}}(\lambda(t_{\mathrm{eq}}), C(t_{\mathrm{eq}}), \beta)}{Z_{\mathcal{S}+\mathcal{E}}(\lambda(0), C(0), \beta)}. \quad (31)$$

This equality, which corresponds to [Eq. (23) in [30]], contains no reference to the dynamics of the composite system. All relations derived in this work—independent of coupling strength or dynamical assumptions—reduce to this equilibrium partition function ratio. The standard JE can be viewed as the specific realization of this identity under Hamiltonian evolution of the composite system. For completeness, the detailed derivation linking Eq. (31) to the conventional JE is provided in Appendix, where we recast Jarzynski's original argument using the present notation.

*Validation.* To validate the framework, we examine how the proposed equalities hold beyond the regime where the JE applies. The model construction is inspired by the
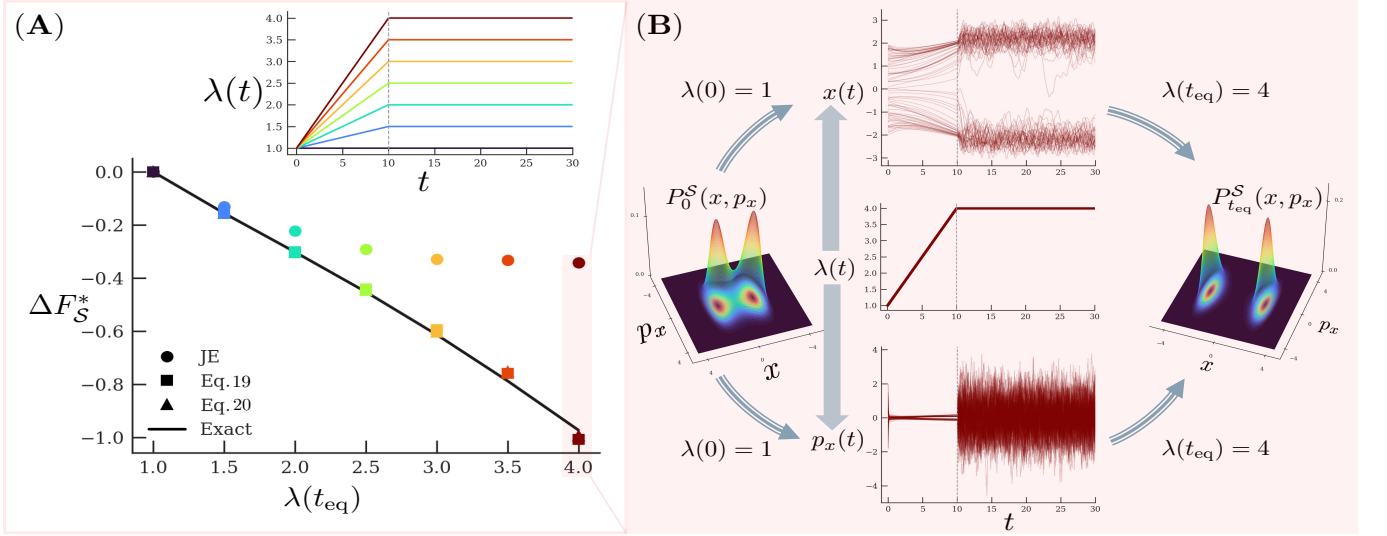
FIG. 2. **(A)** System free energy difference $\Delta F_{\mathcal{S}}^*$ versus the final control value $\lambda(t_{\text{eq}})$: circles (JE), squares (Eq. 19), triangles (Eq. 20), solid line (exact). *Ramp/post-ramp* protocol: $\lambda(t)$ increases linearly from $\lambda(0)$ to the target over $t_\lambda$, then is held fixed. **(B)** Representative run for $\lambda(0) = 1 \rightarrow \lambda(t_{\text{eq}}) = 4$. Center: $\lambda(t)$. Top/bottom: ensemble trajectories of $x(t)$ and $p_x(t)$ (dashed line marks $t_\lambda$). Left/right: canonical marginals $P_0^{\mathcal{S}}(x, p_x)$ and $P_{t_{\text{eq}}}^{\mathcal{S}}(x, p_x)$. During the ramp, dynamics are non-Liouvillian with damping $(\gamma_{\text{nl}}^{\mathcal{S}}, \gamma_{\text{nl}}^{\mathcal{E}})$, post-ramp relaxation is governed by the *underdamped Langevin* equations, whose solution converges to the canonical distribution of the composite system, thereby satisfying the asymptotic equilibration postulate, Eq. (15).

two–stage thought experiment introduced by Jarzynski in his reply [30] to Cohen's critique [22]. In his argument, the control parameter is first driven from $\lambda(0)$ to $\lambda(t_\lambda)$ (stage 1) and then held fixed at $\lambda = \lambda(t_\lambda)$ for relaxation (stage 2). Two observers — one stopping at $t_\lambda$ (out of equilibrium) and the other at $t_{\text{eq}}$ (after relaxation) — record the same work, since no work is performed once the driving protocol is fixed, and calculate the same $\Delta F_{\mathcal{S}}^*$, defined as the equilibrium free energy difference between the canonical states associated with $\lambda(0)$ and $\lambda(t_\lambda) = \lambda(t_{\text{eq}})$ via the HMF partition functions. As Jarzynski emphasized [30], "whether or not we choose to include a relaxation stage has no bearing on the validity" of the JE. While this statement is entirely correct, it implicitly depends on a key assumption about the underlying dynamics. To probe this assumption directly, we construct a validation model consisting of a *non–Liouvillian* ramp (stage 1) that breaks phase–space volume preservation, followed by an *underdamped Langevin* relaxation (stage 2) that enforces the required equilibrium endpoint. In this scenario, the JE no longer holds because phase–space contraction during the ramp breaks the Liouvillian structure of the composite dynamics, preventing the equilibrium partition function ratio in Eq. (31) from reproducing the JE. For a detailed discussion of the validation model and protocol, see the Appendix.

*Conclusions:* In this Letter, we establish exact fluctuation relations for open systems. Using the HMF, the free energy difference is expressed via exponential moments of the HMF shift and an endpoint chi-squared divergence (Thm. 1). The trajectory counterparts of

that relation is obtained under the assumption of asymptotic equilibration (Thm. 2). In the frozen coupling regime the relations reduce to bare-system expressions (Cor. 1) with a clear heat–work decomposition (Cor. 2), and they include the JE and its feedback extensions (GJE) as limiting cases. Validation on a composite model with a non-Liouvillian ramp followed by underdamped Langevin relaxation shows that our endpoint and trajectory relations reproduce the exact free energy differences, demonstrating accuracy and practical applicability beyond the JE regime. The generality of this framework invites direct application to complex nonequilibrium settings where traditional FTs are fragile or inapplicable, including active matter [51, 52], biological processes [53], complex molecular and chemical environments [54], and open quantum systems [55]. In these regimes, our equalities provide a practical route to reconstruct free energy landscapes from realistic driving protocols without imposing idealized dynamical models or weak coupling assumptions. When applied to situations in which the JE is valid, our relations reduce to an exact re-expression of the JE that makes explicit the roles of endpoint overlap and system–environment interaction through the chi-squared factor and the HMF structure. In this sense, the JE becomes a diagnostic tool: deviations can be traced to insufficient overlap [56], rare-event sampling issues [57], or neglected coupling effects [50]. These diagnostics offer concrete guidance for designing and refining experimental and simulation protocols [58], and we expect them to support more reliable free energy estimation across a broad class of strongly driven open systems.

## ACKNOWLEDGMENTS

[1] C. Jarzynski, Nonequilibrium equality for free energy differences, Phys. Rev. Lett. **78**, 2690 (1997).

[2] M. Esposito, U. Harbola, and S. Mukamel, Nonequilibrium fluctuations, fluctuation theorems, and counting statistics in quantum systems, Rev. Mod. Phys. **81**, 1665 (2009).

[3] G. E. Crooks, Nonequilibrium measurements of free energy differences for microscopically reversible Markovian systems, J. Stat. Phys. **90**, 1481 (1998).

[4] G. E. Crooks, Entropy production fluctuation theorem and the nonequilibrium work relation for free energy differences, Phys. Rev. E **60**, 2721 (1999).

[5] J. Liphardt, S. Dumont, S. B. Smith, I. Tinoco Jr, and C. Bustamante, Equilibrium information from nonequilibrium measurements in an experimental test of Jarzynski's equality, Science **296**, 1832 (2002).

[6] M. Ohzeki, Quantum annealing with the Jarzynski equality, Phys. Rev. Lett. **105**, 050401 (2010).

[7] T. Xiong, L. Yan, F. Zhou, K. Rehan, D. Liang, L. Chen, W. Yang, Z. Ma, M. Feng, and V. Vedral, Experimental verification of a Jarzynski-related information-theoretic equality by a single trapped ion, Phys. Rev. Lett. **120**, 010601 (2018).

[8] W. Liu, Z. Niu, W. Cheng, X. Li, C.-K. Duan, Z. Yin, X. Rong, and J. Du, Experimental test of the Jarzynski equality in a single spin-1 system using high-fidelity single-shot readouts, Phys. Rev. Lett. **131**, 220401 (2023).

[9] N. C. Harris, Y. Song, and C.-H. Kiang, Experimental free energy surface reconstruction from single-molecule force spectroscopy using Jarzynski's equality, Phys. Rev. Lett. **99**, 068101 (2007).

[10] F. Douarche, S. Ciliberto, A. Petrosyan, and I. Rabbiosi, An experimental test of the Jarzynski equality in a mechanical experiment, Europhys. Lett. **70**, 593 (2005).

[11] O.-P. Saira, Y. Yoon, T. Tanttu, M. Möttönen, D. Averin, and J. P. Pekola, Test of the Jarzynski and Crooks fluctuation relations in an electronic system, Phys. Rev. Lett. **109**, 180601 (2012).

[12] G. Wimsatt, O.-P. Saira, A. B. Boyd, M. H. Matheny, S. Han, M. L. Roukes, and J. P. Crutchfield, Harnessing fluctuations in thermodynamic computing via time-reversal symmetries, Phys. Rev. Researh **3**, 033115 (2021).

[13] J. Gore, F. Ritort, and C. Bustamante, Bias and error in estimates of equilibrium free-energy differences from nonequilibrium measurements, Proc. Natl. Acad. Sci. U.S.A. **100**, 12564 (2003).

[14] S. Toyabe, T. Sagawa, M. Ueda, E. Muneyuki, and M. Sano, Experimental demonstration of information-to-energy conversion and validation of the generalized Jarzynski equality, Nat. Phys. **6**, 988 (2010).

[15] J. Deng, J. D. Jaramillo, P. Hänggi, and J. Gong, Deformed Jarzynski equality, Entropy **19**, 419 (2017).

[16] S. Hernández-Gómez, S. Gherardini, A. Belenchia, A. Trombettoni, M. Paternostro, and N. Fabbri, Experimental signature of initial quantum coherence on entropy production, npj Quantum Information **9**, 86 (2023).

[17] D. Hahn, M. Dupont, M. Schmitt, D. J. Luitz, and M. Bukov, Quantum many-body Jarzynski equality and dissipative noise on a digital quantum computer, Phys. Rev. X. **13**, 041023 (2023).

[18] S. Pressé and R. Silbey, Ordering of limits in the Jarzynski equality, J. Chem. Phys. **124** (2006).

[19] L. Chen, On the Crooks fluctuation theorem and the Jarzynski equality, J. Chem. Phys. **129** (2008).

[20] G. E. Crooks, Comment regarding" on the Crooks fluctuation theorem and the Jarzynski equality"[j. chem. phys. 129, 091101 (2008)] and" nonequilibrium fluctuation-dissipation theorem of Brownian dynamics"[j. chem. phys. 129, 144113 (2008)], (2009).

[21] T. Sagawa and M. Ueda, Generalized Jarzynski equality under nonequilibrium feedback control, Phys. Rev. Lett. **104**, 090602 (2010).

[22] E. Cohen and D. Mauzerall, A note on the Jarzynski equality, J. Stat. Mech.: Theory Exp. **2004** (07), P07006.

[23] E. Cohen and D. Mauzerall, The Jarzynski equality and the Boltzmann factor, Mol. Phys. **103**, 2923 (2005).

[24] A. Argun, A.-R. Moradi, E. Pince, G. B. Bagci, and G. Volpe, Experimental evidence of the failure of Jarzynski equality in active baths, arXiv preprint arXiv:1601.01123 (2016).

[25] J. M. Vilar and J. M. Rubi, Failure of the work-Hamiltonian connection for free-energy calculations, Phys. Rev. Lett. **100**, 020601 (2008).

[26] B. Palmieri and D. Ronis, Jarzynski equality: Connections to thermodynamics and the second law, Phys. Rev. E **75**, 011133 (2007).

[27] A. M. Monge, M. Manosas, and F. Ritort, Experimental test of ensemble inequivalence and the fluctuation theorem in the force ensemble in DNA pulling experiments, Phys. Rev. E **98**, 032146 (2018).

[28] A. Sone, Y.-X. Liu, and P. Cappellaro, Quantum Jarzynski equality in open quantum systems from the one-time measurement scheme, Phys. Rev. Lett. **125**, 060602 (2020).

[29] C. Jarzynski, Equilibrium free-energy differences from nonequilibrium measurements: A master-equation approach, Phys. Rev. E **56**, 5018 (1997).

[30] C. Jarzynski, Nonequilibrium work theorem for a system strongly coupled to a thermal environment, J. Stat. Mech.: Theory Exp. **2004** (09), P09005.

[31] T. Speck and U. Seifert, The Jarzynski relation, fluctuation theorems, and stochastic thermodynamics for non-Markovian processes, J. Stat. Mech.: Theory Exp. **2007** (09), L09002.

[32] K. Maruyama, F. Nori, and V. Vedral, Colloquium: The physics of Maxwell's demon and information, Rev. Mod.

Phys. **81**, 1 (2009).

[33] V. Serreli, C.-F. Lee, E. R. Kay, and D. A. Leigh, A molecular information ratchet, Nature **445**, 523 (2007).

[34] C. Maes, Local detailed balance, SciPost Phys. Lect. Notes. , 032 (2021).

[35] E. Schöll-Paschinger and C. Dellago, A proof of Jarzynski's nonequilibrium work theorem for dynamical systems that conserve the canonical distribution, J. Chem. Phys. **125** (2006).

[36] P. Gaspard and R. Kapral, Active matter, microreversibility, and thermodynamics, Research (2020).

[37] I. A. Martínez, G. Bisker, J. M. Horowitz, and J. M. Parrondo, Inferring broken detailed balance in the absence of observable currents, Nat. Commun. **10**, 3542 (2019).

[38] F. S. Gnesotto, F. Mura, J. Gladrow, and C. P. Broedersz, Broken detailed balance and non-equilibrium dynamics in living systems: a review, Rep. Prog. Phys. **81**, 066601 (2018).

[39] K. Goswami, Work fluctuation relations for a dragged Brownian particle in active bath, Physica **525**, 223 (2019).

[40] A. Godec and D. E. Makarov, Challenges in inferring the directionality of active molecular processes from single-molecule fluorescence resonance energy transfer trajectories, J. Phys. Chem. Lett. **14**, 49 (2022).

[41] S. Krishnamurthy, S. Ghosh, D. Chatterji, R. Ganapathy, and A. Sood, A micrometre-sized heat engine operating between bacterial reservoirs, Nat. Phys. **12**, 1134 (2016).

[42] P. Reimann, Brownian motors: noisy transport far from equilibrium, Phys. Rep. **361**, 57 (2002).

[43] C. Bechinger, R. Di Leonardo, H. Löwen, C. Reichhardt, G. Volpe, and G. Volpe, Active particles in complex and crowded environments, Rev. Mod. Phys. **88**, 045006 (2016).

[44] L. Caprini, U. Marini Bettolo Marconi, and A. Puglisi, Activity induced delocalization and freezing in self-propelled systems, Sci. Rep. **9**, 1386 (2019).

[45] T. Vicsek and A. Zafeiris, Collective motion, Phys. Rep. **517**, 71 (2012).

[46] M. E. Cates and J. Tailleur, Motility-induced phase separation, Annu. Rev. Condens. Matter Phys. **6**, 219 (2015).

[47] É. Fodor, C. Nardini, M. E. Cates, J. Tailleur, P. Visco, and F. Van Wijland, How far from equilibrium is active matter?, Phys. Rev. Lett. **117**, 038103 (2016).

[48] A. Argun, A.-R. Moradi, E. Pinçe, G. B. Bagci, A. Imparato, and G. Volpe, Non-Boltzmann stationary distributions and nonequilibrium relations in active baths, Phys. Rev. E. **94**, 062150 (2016).

[49] M. Rahbar and C. J. Stein, Thermodynamic potentials from a probabilistic view on the system-environment interaction energy, arXiv preprint arXiv:2505.00188 (2025).

[50] P. Talkner and P. Hänggi, Colloquium: Statistical mechanics and thermodynamics at strong coupling: Quantum and classical, Rev. Mod. Phys. **92**, 041002 (2020).

[51] M. K. Johnsrud and R. Golestanian, Fluctuation dissipation relations for active field theories, Phys. Rev. Res. **7**, L032053 (2025).

[52] R. Bebon, J. F. Robinson, and T. Speck, Thermodynamics of active matter: Tracking dissipation across scales, Phys. Rev. X **15**, 021050 (2025).

[53] M. T. Woodside and S. M. Block, Reconstructing folding energy landscapes by single-molecule force spectroscopy, Annu. Rev. Biophys. **43**, 19 (2014).

[54] M. P. Leighton and D. A. Sivak, Flow of energy and information in molecular machines, Annu. Rev. Phys. Chem **76** (2025).

[55] K. Beyer and W. T. Strunz, Operational work fluctuation theorem for open quantum systems, Phys. Rev. Lett. **134**, 140403 (2025).

[56] L. Schmidt, C. J. Wilson, S. Behera, and B. L. de Groot, Free energy calculations for protein design, (2025).

[57] Z. Kuang, K. M. Singh, D. J. Oliver, P. B. Dennis, C. C. Perry, and R. R. Naik, Gamma estimator of jarzynski equality for recovering binding energies from noisy dynamic data sets, Nat. Commun. **11**, 5517 (2020).

[58] C. Casert and S. Whitelam, Learning protocols for the fast and efficient control of active matter, Nat. Commun. **15**, 9128 (2024).

[59] S. Whitelam, Improving noisy free-energy measurements by adding more noise, Phys. Rev. E **112**, 014133 (2025).

[60] P. Hänggi, P. Talkner, and M. Borkovec, Reaction-rate theory: fifty years after kramers, Reviews of modern physics **62**, 251 (1990).

[61] S. Kieninger and B. G. Keller, GROMACS stochastic dynamics and BAOAB are equivalent configurational sampling algorithms, J. Chem. Theory Comput. **18**, 5792 (2022).

# End Matter

*Appendix:* We collect detailed proofs of the theorems and corollaries stated in the main text. Each result is reproduced in full mathematical form to ensure completeness and clarity of the derivations.

*Proof of Theorem 1.* From Eqs. (5), (6) and (10), we have

$$e^{-\beta \Delta \mathcal{H}_\beta^*(X_\mathcal{S})} = \frac{P_{t_{\text{eq}}}^\mathcal{S}(X_\mathcal{S}, \beta)\, \mathcal{Z}_\mathcal{S}^*(\lambda(t_{\text{eq}}), C(t_{\text{eq}}), \beta)}{P_0^\mathcal{S}(X_\mathcal{S}, \beta)\, \mathcal{Z}_\mathcal{S}^*(\lambda(0), C(0), \beta)}. \quad (32)$$

Averaging Eq. (32) over the final ensemble, Eq. (11), yields

$$\left\langle e^{-\beta \Delta \mathcal{H}_\beta^*} \right\rangle_\mathcal{S}$$
$$= \frac{\mathcal{Z}_\mathcal{S}^*(\lambda(t_{\text{eq}}), C(t_{\text{eq}}), \beta)}{\mathcal{Z}_\mathcal{S}^*(\lambda(0), C(0), \beta)} \int \mathrm{d}X_\mathcal{S}\, \frac{\left(P_{t_{\text{eq}}}^\mathcal{S}(X_\mathcal{S}, \beta)\right)^2}{P_0^\mathcal{S}(X_\mathcal{S}, \beta)}. \quad (33)$$

Using Eq. (9) for the partition function ratio and Eq. (12) for the divergence, Eq. (33) reduces to

$$\left\langle e^{-\beta \Delta \mathcal{H}_\beta^*} \right\rangle_\mathcal{S} = e^{-\beta \Delta F_\mathcal{S}^*(\beta)} \left[1 + \chi^2\left(P_{t_{\text{eq}}}^\mathcal{S} \parallel P_0^\mathcal{S}\right)\right], \quad (34)$$

which rearranges to Eq. (7). For the positive exponential, Eqs. (5), (6) and (10), give

$$e^{+\beta \Delta \mathcal{H}_\beta^*(X_\mathcal{S})} = \frac{P_0^\mathcal{S}(X_\mathcal{S}, \beta)\, \mathcal{Z}_\mathcal{S}^*(\lambda(0), C(0), \beta)}{P_{t_{\text{eq}}}^\mathcal{S}(X_\mathcal{S}, \beta)\, \mathcal{Z}_\mathcal{S}^*(\lambda(t_{\text{eq}}), C(t_{\text{eq}}), \beta)}. \quad (35)$$

Averaging Eq. (35) with Eq. (11) cancels the denominator and yields

$$\left\langle e^{+\beta \Delta \mathcal{H}_\beta^*} \right\rangle_\mathcal{S} = \frac{\mathcal{Z}_\mathcal{S}^*(\lambda(0), C(0), \beta)}{\mathcal{Z}_\mathcal{S}^*(\lambda(t_{\text{eq}}), C(t_{\text{eq}}), \beta)} = e^{+\beta \Delta F_\mathcal{S}^*(\beta)}, \quad (36)$$

which reproduces Eq. (8).  □

*Theorem 2: (proof of Eq. (14)).* Starting from Eq. (13),

$$P_t(X, \beta) = \int_{\Gamma_0} \mathrm{d}X_0\, \delta\big(X - \mathcal{T}_t(X_0)\big)\, P_0(X_0, \beta), \quad (37)$$

where $X = (X_\mathcal{S}, X_\mathcal{E})$, the system marginal is obtained by integrating over $X_\mathcal{E}$. After exchanging the order of integration and applying the sifting property, we obtain

$$P_t^\mathcal{S}(X_\mathcal{S}, \beta) = \int \mathrm{d}X_\mathcal{E}\, P_t(X_\mathcal{S}, X_\mathcal{E}, \beta) = \int \mathrm{d}X_\mathcal{E} \int_{\Gamma_0} \mathrm{d}X_0\, \delta\big(X_\mathcal{S} - \mathcal{T}_t^\mathcal{S}(X_0)\big) \delta\big(X_\mathcal{E} - \mathcal{T}_t^\mathcal{E}(X_0)\big)\, P_0(X_0, \beta)$$

$$= \int_{\Gamma_0} \mathrm{d}X_0\, P_0(X_0, \beta) \int \mathrm{d}X_\mathcal{E}\, \delta\big(X_\mathcal{S} - \mathcal{T}_t^\mathcal{S}(X_0)\big) \delta\big(X_\mathcal{E} - \mathcal{T}_t^\mathcal{E}(X_0)\big) = \int_{\Gamma_0} \mathrm{d}X_0\, P_0(X_0, \beta)\, \delta\big(X_\mathcal{S} - \mathcal{T}_t^\mathcal{S}(X_0)\big),$$

which reproduces Eq. (14).  □

*Theorem 2: (proof of numerators in Eqs. 16 and 17).* Starting from Eq. (15) and using Eq. (11) and Eq. (14), we can write

$$\left\langle e^{-\beta \Delta \mathcal{H}_\beta^*(X_\mathcal{S})} \right\rangle_\mathcal{S} = \lim_{t \to t_{\text{eq}}} \int \mathrm{d}X_\mathcal{S}\, e^{-\beta \Delta \mathcal{H}_\beta^*(X_\mathcal{S})}\, P_t^\mathcal{S}(X_\mathcal{S}, \beta) = \int_{\Gamma_0} \mathrm{d}X_0\, e^{-\beta \Delta \mathcal{H}_\beta^*\left(\mathcal{T}_{t_{\text{eq}}}^\mathcal{S}(X_0)\right)}\, P_0(X_0, \beta)$$

$$= \left\langle e^{-\beta \Delta \mathcal{H}_\beta^*\left(\mathcal{T}_{t_{\text{eq}}}^\mathcal{S}(X_0)\right)} \right\rangle_{X_0}, \quad (38)$$

which confirms the numerator in Eq. (16). Using the same strategy, we obtain the numerator in Eq. (17) .  □

*JE as limiting case.* Starting from the Eq. (31), and assuming the composite system evolves under Hamilton's equations generated by $\mathcal{H}_{\mathcal{S}+\mathcal{E}}$. The instantaneous energy change along a trajectory $X(t|X_0)$ obeys

$$\frac{\mathrm{d}}{\mathrm{d}t}\, \mathcal{H}_{\mathcal{S}+\mathcal{E}}(X(t|X_0), \lambda(t), C(t))$$
$$= \underbrace{\nabla_X \mathcal{H}_{\mathcal{S}+\mathcal{E}} \cdot \dot{X}}_{\text{advection}} + \frac{\partial \mathcal{H}_\mathcal{S}}{\partial \lambda}\, \dot{\lambda} + \frac{\partial \mathcal{V}_{\mathcal{S}\mathcal{E}}}{\partial C}\, \dot{C}. \quad (39)$$

The advection term vanishes because the Hamiltonian flow is symplectic, $\nabla_X \mathcal{H}_{\mathcal{S}+\mathcal{E}} \cdot \dot{X} = \{\mathcal{H}_{\mathcal{S}+\mathcal{E}}, \mathcal{H}_{\mathcal{S}+\mathcal{E}}\} = 0$. Integrating over $t \in [0, t_{\text{eq}}]$ gives

$$\mathcal{H}_{\mathcal{S}+\mathcal{E}}(X(t|X_0), \lambda(t), C(t)) - \mathcal{H}_{\mathcal{S}+\mathcal{E}}(X(0|X_0), \lambda(0), C(0))$$
$$= \int_0^{t_{\text{eq}}} \mathrm{d}t \left(\frac{\partial \mathcal{H}_\mathcal{S}}{\partial \lambda}\, \dot{\lambda} + \frac{\partial \mathcal{V}_{\mathcal{S}\mathcal{E}}}{\partial C}\, \dot{C}\right) =: W_\lambda + W_C. \quad (40)$$

Using the above equation and Liouville's theorem, which guarantees a unit Jacobian for the change of variables $(X_0 \mapsto X)$, Jarzynski demonstrated that the equilibrium

partition function ratio for the composite system can be expressed as

$$\left\langle e^{-\beta(W_\lambda + W_C)}\right\rangle_{X_0,\beta} = \frac{Z_{\mathcal{S}+\mathcal{E}}(\lambda(t_{\text{eq}}), C(t_{\text{eq}}), \beta)}{Z_{\mathcal{S}+\mathcal{E}}(\lambda(0), C(0), \beta)}. \quad (41)$$

For frozen coupling $C(t) \equiv C(0)$, Eq. (41) reduces to

$$e^{-\beta\Delta F_{\mathcal{S}}^*} = \left\langle e^{-\beta W_\lambda}\right\rangle_{X_0,\beta}, \quad (42)$$

which is the conventional JE expressed in the HMF notation. □

*Validation model and protocol*: *Model specification.* We validate Eqs. (19) and (20) on an analytically tractable composite system consisting of a double well potential [59, 60] coupled to a harmonic environment. The system Hamiltonian is $\mathcal{H}_{\mathcal{S}}(x, p_x; \lambda) = \frac{1}{2m}p_x^2 + U_{\mathcal{S}}(x; \lambda)$, where $U_{\mathcal{S}}(x; \lambda) = \frac{1}{4}(x^2 - \lambda)^2$, the environment Hamiltonian $\mathcal{H}_{\mathcal{E}}(y, p_y) = \frac{1}{2m}p_y^2 + \frac{1}{2}\omega^2 y^2$, and the interaction $\mathcal{V}_{\mathcal{S}\mathcal{E}}(x, y; C) = Cxy$ with coupling constant $C$. The total Hamiltonian follows Def. 2. Eq. (4) gives $\mathcal{H}_{\beta}^*(x, p_x; \lambda, C)$, the corresponding marginal and partition function follow from $P^{\mathcal{S}}(X_{\mathcal{S}}, \beta) = e^{-\beta\mathcal{H}_{\beta}^*(X_{\mathcal{S}}, \lambda, C)}/\mathcal{Z}_{\mathcal{S}}^*(\lambda, C, \beta)$ and $\mathcal{Z}_{\mathcal{S}}^*(\lambda, C, \beta) = \int dX_{\mathcal{S}} e^{-\beta\mathcal{H}_{\beta}^*(X_{\mathcal{S}}, \lambda, C)}$. The composite system is initialized in the canonical ensemble at $(\lambda(0), C(0))$, ensuring exact sampling from Eq. (2). *Two–stage protocol.* We work in the frozen–coupling regime $C(t) \equiv C(0)$ and vary only $\lambda(t)$. During the ramp $0 \leq t \leq t_\lambda$, $\lambda(t)$ changes linearly from $\lambda(0)$ to $\lambda(t_\lambda)$ (see Fig. 2). The composite evolves under non-Liouvillian deterministic dynamics $\dot{x} = p_x$, $\dot{p}_x = -\partial_x U_{\mathcal{S}}(x; \lambda(t)) - Cy - \gamma_{\text{nl}}^{\mathcal{S}} p_x$, $\dot{y} = p_y$, $\dot{p}_y = -\omega^2 y - Cx - \gamma_{\text{nl}}^{\mathcal{E}} p_y$, where $\gamma_{\text{nl}}^{\mathcal{S}}$ and $\gamma_{\text{nl}}^{\mathcal{E}}$ denote constant damping coefficients used only during the ramp. At $t = t_\lambda$, $\lambda(t)$ is held fixed and the composite relaxes under underdamped Langevin dynamics [60] obeying the FDT: $\dot{x} = p_x$, $\dot{p}_x = -\partial_x U_{\mathcal{S}}(x; \lambda(t_{\text{eq}})) - Cy - \gamma p_x + \eta_x(t)$, $\dot{y} = p_y$, $\dot{p}_y = -\omega^2 y - Cx - \gamma p_y + \eta_y(t)$, where $\gamma$ is the friction coefficient. The terms $\eta_x(t)$ and $\eta_y(t)$ are independent standard Gaussian white noises with, $\langle \eta_i(t) \rangle = 0$ and $\langle \eta_i(t)\eta_j(t') \rangle = (2\gamma/\beta)\delta_{ij}\delta(t - t')$ for $i, j \in \{x, y\}$. These equations converge to the canonical distribution associated with the frozen Hamiltonian $\mathcal{H}_{\mathcal{S}+\mathcal{E}}(x, p_x, y, p_y; \lambda(t_{\text{eq}}), C)$, satisfying the asymptotic equilibration postulate Eq. (15). *Estimators.* After relaxation, the endpoint relations Eqs. (19) and (20) are evaluated. The "Exact" curve is the equilibrium free energy difference from Eq. (30). For the JE estimator we accumulate the ramp work

$$W_\lambda = \int_0^{t_\lambda} \partial_\lambda U_{\mathcal{S}}(x_t; \lambda_t)\,\dot{\lambda}\,dt, \quad (43)$$

with frozen $C$, and compute, $e^{-\beta\Delta F_{\mathcal{S}}^*} \approx \langle e^{-\beta W_\lambda} \rangle$. Fig. 2 reports $\Delta F_{\mathcal{S}}^*$ as a function of $\lambda(t_{\text{eq}})$: the JE estimator deviates because the Liouvillian assumption is violated, whereas the endpoint equalities reproduce the exact equilibrium free energies. As a consistency check, when we set $\gamma_{\text{nl}}^{\mathcal{S}} = \gamma_{\text{nl}}^{\mathcal{E}} = 0$ (Hamiltonian/Liouvillian ramp), the JE work estimator, the endpoint equalities, and the equilibrium reference all coincide, recovering the standard JE (see Fig.3). For completeness, all numerical parameters used for the simulations are listed in Table. I.

TABLE I. Numerical inputs for the validation runs (frozen coupling; non-Liouvillian ramp, under damped Langevin relaxation). All quantities are in reduced, dimensionless units.

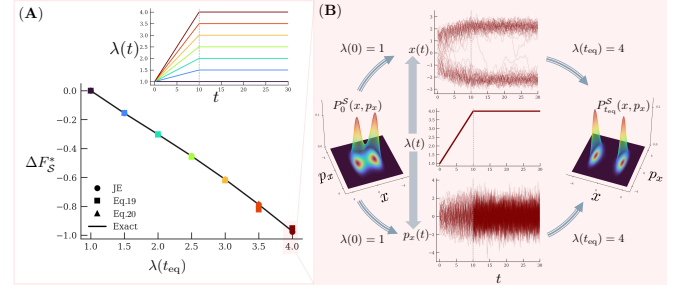| Quantity | Value |
| --- | --- |
| $m$ (Mass) | 1.0 |
| System–environment coupling $C$ | 1.0 |
| $\beta$ | 1.0 |
| System–environment coupling $C$ | 1.0 |
| Environment frequency $\omega$ | 1.0 |
| Initial protocol $\lambda(0)$ | 1.0 |
| Ramp duration $t_{\text{ramp}}$ | 10.0 |
| Relaxation duration $t_{\text{eq}} - t_\lambda$ | 20.0 |
| Ramp time step $\Delta t$ | $2 \times 10^{-3}$ |
| Relaxation time step $\Delta t_{\text{relax}}$ | $2 \times 10^{-2}$ |
| Non–Liouvillian drag (system) $\gamma_{\text{nl}}^{\mathcal{S}}$ | 20.0 |
| Non–Liouvillian drag (environment) $\gamma_{\text{nl}}^{\mathcal{E}}$ | 5.0 |
| Langevin damping (relaxation) $\gamma$ | 2.0 |
| Initial ensemble size $N_{\text{traj}}$ | 30000 |
| Ramp integrator | RK4 |
| Relax integrator | BAOAB [61] |



FIG. 3. Same plotting conventions as in Fig. 2. During the ramp we set $\gamma_{\text{nl}}^{\mathcal{S}} = \gamma_{\text{nl}}^{\mathcal{E}} = 0$ (Hamiltonian, Liouvillian, volume-preserving).