# Non-Markovian nonequilibrium information dynamics

Qian Zeng[1] and Jin Wang[1,2,*]

[1]*State Key Laboratory of Electroanalytical Chemistry, Changchun Institute of Applied Chemistry, Changchun, Jilin 130022, China*
[2]*Department of Chemistry and Physics, State University of New York, Stony Brook, New York 11794, USA*

We probe into the dynamics of interacting non-Markovian information systems. The stochastic dynamics of information has two aspects: the self-evolution and interaction. We show that self-evolution of a non-Markovian information system can be described by a Markov-type master equation with memory dependence. We also reveal that the interaction between systems can be fully embodied into the information dynamics of the composite information system. To characterize time irreversibility of the self-evolution and the interaction, we apply the landscape-flux theory to both stochastic and thermal information dynamics. The driving force of the nonequilibrium information dynamics can be decomposed into time-reversible (detailed balance preserving landscape part) and -irreversible (detailed balance breaking nonequilibrium flux part) parts. The time-irreversible part of the driving force fully depicts the time-irreversibility behavior in the stochastic dynamics. The time irreversibility of the interactions between systems reflected in nonequilibrium thermodynamics can be seen in the decomposition of the mutual information rate which corresponds to decomposition of the driving force. In particular, the time-irreversible part of mutual information rate reveals the underlying relationship among the entropy production rates of the information systems. We propose the finite memory approximation method and demonstrate that the above mentioned features can be found in a wide class of non-Markovian nonequilibrium information systems. Finally, we derive the lower and upper bounds for informational entities under concern with clear meanings.

## I. INTRODUCTION

Studies on the nonequilibrium behaviors of two interacting systems with finite states have shown their importance in mesoscopic and microscopic information dynamics [1–6]. Usually, the dynamics of interacting systems in random environments with time-invariant parameters (temperatures, chemical potentials, etc.) and infinite degrees of freedom is always considered to be non-Markovian with finite memories [7]. Although the usual analytical and numerical methods for Markov processes can be also applied for getting the pictures of both stochastic dynamics and thermodynamics of a composite Markov system, it has been proven to be difficult that we can depict the behaviors of the subsystems with the same ingredient. This is because the two subsystems may also be random environments with time-variant parameters for each other due to the comparable sizes and state-switching rates of the subsystems. This indicates that none of the subsystems need to be Markovian. Because of the existence of the interactions, every subsystem has to adjust itself to adapt to the time-variant environment (the other subsystem) and then the corresponding process shows a remarkable path dependence by summing away the degree of freedom of the other subsystem, which means it has a memory. The information of the past states in a memory are always embedded into the parameters of the dynamics of the system to determine

the present state, where the parameters are therefore time variant. This memory is also always finite because historical information of the systems is dissipated into the random environments. This always leads to complicated noncommutative matrices multiplications in mathematical treatments [8]. Due to the difficulties on dealing with noncommutative matrices multiplications [9], there is no analytical method to reveal the underlying physical mechanisms which give rise to some peculiar time-irreversible or nonequilibrium behaviors of such interacting systems.

A recent progress in study of bivariate Markov chains by Zeng and Wang [10] has shown that if both the subsystems behave as Markov chains, the stochastic forces behind the two subsystems can be divided into two parts. One part can be regarded as a gradientlike force that attracts a subsystem down to each state and due to potential landscape according to the steady state distribution of the states. The other part is a curling force that drives the subsystem to rotate among states with steady probability fluxes. By restricting the subsystems to be Markov systems, this decomposition of the forces which is the so-called "landscape-flux" decomposition [11–13] is embodied in decomposing the transition matrices of the subsystems into their time-reversible and -irreversible parts. Although this model is restrictive, the landscape-flux decomposition sets up a formal theory to deal with more general cases where the subsystems and the composite system could be non-Markovian.

In this paper, we probe into stochastic dynamics and thermodynamics of non-Markovian interacting information systems where for no loss of generality we assume that the

---

*To whom correspondence should be addressed: jin.wang.1@stonybrook.edu

two information systems and the composite system have a unified memory length of $m$. In stochastic dynamics, we slice sequences of a system into time-successive memories of length $m$. By regarding all possible memories as generalized states, we construct the Markov-type master equation which characterizes the evolution of distribution of memories. The force behind the information dynamics is considered to be the transition probability that a memory "jumps" to another time-successive memory. Based on this, we define the probability flux of a non-Markovian information system which fully characterizes the time irreversibility in stochastic dynamics. Geometrically, this probability flux can be interpreted as an tensor of order $m + 1$ in a linear tensor space. It carries informational observables corresponding to the memories and flows from the system to the nonequilibrium environments. On the other hand, the negative logarithmic distribution of the memories in steady state is the so-called potential landscape which preserves some informational observables against the nonequilibrium behaviors, i.e., the probability flux. Correspondingly, we divide the driving force (transition probability) into the time-reversible and -irreversible parts to uncover the underlying landscape-flux decomposition of the non-Markovian dynamics. The interactions between the two systems are fully embodied into the dynamics of the composite system which may also be non-Markovian with memory. The probability fluxes of the subsystems can be regarded as the coarse-grained version of that of the composite system. This builds up the bridge between the time irreversibility of the subsystems and the composite system.

To characterize the interactions of the information systems, we focus on the mutual information rate (MIR) [14,15] which depicts the information of interaction between subsystems. The explicit form of MIR can be obtained from the master equations of systems. It is the average over the detailed interactions with respect to successive memory transitions. Characterization of time irreversibility in thermodynamics follows the spirit of landscape-flux decomposition in stochastic dynamics. The MIR can be divided into time-reversible and -irreversible parts explicitly corresponding to the decomposition of stochastic forces. It can be seen from the stochastic thermodynamics [16] that the time-irreversible part of the MIR which is driven by the probability flux is the increasing or decreasing rate of the interactions between subsystems in nonequilibrium environments. On the other hand, the time-irreversible part of the MIR can be viewed as the preserved correlations between subsystems.

In thermodynamics, the time irreversibility or nonequilibriumness of a system is measured by the entropy production rate (EPR) which can be used to quantify the rate of information dissipation from a system into the nonequilibrium environments. Influences on the EPR of one subsystem through the environments and the other subsystem can be seen from the time-irreversible part of MIR. Then, it can be concluded that the time-irreversible part of MIR quantifies the rate of interacting information dissipation of systems under nonequilibrium condition.

The explicit forms of the physical entities in both stochastic dynamics and thermodynamics presented in this paper allow us to evaluate these entities easily in numerical analysis. With these explicit expressions for non-Markovian cases, we

can even analyze the dynamical behavior of the interacting information systems with unknown memory lengths, provided the dynamics of the composite system. This is the so-called "finite memory approximation" (FMA) which is based on the martingale convergence theorem [17] for stationary and ergodic processes. Theoretically, FMA can approximate these physical entities of a non-Markov system in thermodynamics with arbitrary precisions.

In some situations, we may not obtain adequate information of the systems. It then becomes unworthy (or even impossible) to evaluate the exact values of the system observables under concern. Then, it would be meaningful to check whether the designed systems have the desired properties of entities within certain bounds rather than to evaluate the exact values for practical applications. For this reason, we derive the upper and lower bounds of the EPRs and the time-irreversible MIR by using the log-sum inequality.

To demonstrate the power of the analysis (explicit expressions and FMA) for non-Markovian cases in this paper, we construct an example of two interacting non-Markovian information systems, which involves the information dissipation and (feedback) control. Since this model can be solved analytically, we clarify the meanings of the above non-Markovian and nonequilibrium entities by detailed discussions on the dynamical observables of the systems: the optimal code lengths of the system states.

## II. STOCHASTIC DYNAMICS OF INTERACTING NON-MARKOVIAN SYSTEMS

### A. Formulation of self-evolution, time-irreversibility characterization, and landscape-flux decomposition

Consider an open system $X$ coupled to random environments and generates a finite-state, discrete-time, and irreducible chain with state space $\mathcal{X} = \{1, 2, \dots, l\}$: we say that it has memories with length $m$ because the conditional probability $q_x(x_t | x_1, x_2, \dots, x_{t-1}) = q_x(x_t | x_{t-m}, x_{t-m+1}, \dots, x_{t-1})$ for arbitrary state $x \in \mathcal{X}$ and $t > m$. This means the state $x_t$ only depends on the sequence $[x_{t-m}, x_{t-m+1}, \dots, x_{t-1}]$ with length up to $m$. This sequence is the so-called memory of $X$. Here, subscripts of states $x$ represent occurrence orders of $x$ in time. Since historical information of an open system is dissipated into environments with infinite degrees of freedom, we mainly focus on the dynamics of $X$ to be non-Markovian with finite $m$ ($m > 1$).

We note that the Markov chains can be be regarded as a special case in this paper while the memory length $m$ equals to 1. Then, the corresponding conditional probability $q_x(x_t | x_{t-m}, x_{t-m+1}, \dots, x_{t-1})$ reduces to $q_x(x_t | x_{t-1})$ which is the so-called transition probability of the Markov chain. The transition probability $q_x(x_t | x_{t-1})$ quantifies the probability of the transition between state $x_{t-1}$ and state $x_t$ representing the underlying stochastic dynamics of a Markov chain. When the system is non-Markovian, i.e., the memory length $m > 1$, the conditional probability $q_x(x_t | x_{t-m}, x_{t-m+1}, \dots, x_{t-1})$ quantifies the probability of the transition between a state sequence with the memory length $m$ $\chi_{t-1} = [x_{t-m}, x_{t-m+1}, \dots, x_{t-1}]$ and another time-successive state sequence with the memory length
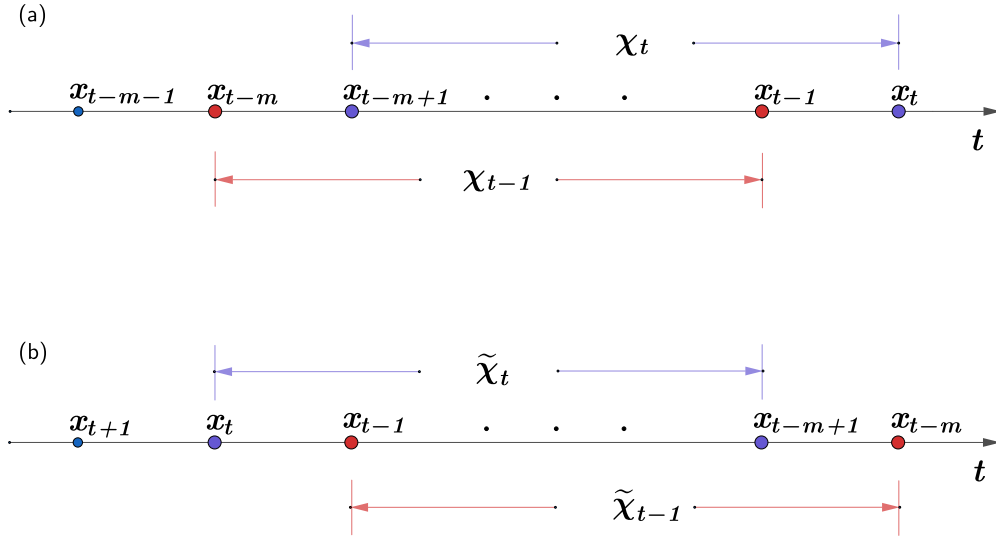
FIG. 1. Locations of two time-successive memories in a sequence of non-Markovian system $X$. (a) Locations of $(\chi_{t-1}, \chi_t)$ along the forward time arrow. (b) Locations of the time reversals $(\widetilde{\chi}_t, \widetilde{\chi}_{t-1})$ along the backward time arrow.

$m$  $\chi_t = [x_{t-m+1}, x_{t-m+2}, \ldots, x_t]$ representing the underlying stochastic non-Markovian dynamics. Here, "time successive" means that the two sequences $\chi_{t-1}$ and $\chi_t$ intersect with each other in time such that $\chi_t \cap \chi_{t-1} = [x_{t-m+1}, x_{t-m+2}, \ldots, x_{t-1}]$ (see Fig. 1). Then, it is reasonable to rewrite $q_x$ as $q_x(\chi_t|\chi_{t-1}) \equiv q_x(x_t|x_{t-m}, x_{t-m+1}, \ldots, x_{t-1})$. This $q_x$ can be regarded as the transition probability between time-successive sequences of length $m$ of a non-Markovian system, which is similar to the Markovian case where we can find that $\chi_{t-1} = x_{t-1}$ and $\chi_t = x_t$. Due to the description of the transition probability $q_x$, the characterization of the dynamics of a non-Markovian system should be based on the transitions between state sequences with memory length $m$ (they are considered as the memories of the future) rather than transitions between states. A Markov-type master equation for the underlying non-Markovian evolution of distribution with memories $\chi$ can come out based on this reasoning. It can be expressed as

$$\Pr(\chi_t) = \sum_{\chi_{t-1}} q_x(\chi_t|\chi_{t-1})\Pr(\chi_{t-1}), \qquad (1)$$

where the operator $[q_x(\chi_t|\chi_{t-1})]$ is the stochastic propagator of the $X$ process. The solution of this equation can be given by $\Pr(\chi_t) = \sum_{\chi_1,\ldots,\chi_{t-1}} \prod_{i=m+1}^{t} q_x(\chi_i|\chi_{i-1})\Pr(\chi_m)$, provided the initial distribution $\Pr(\chi_m)$. Due to the complexity of eigenspace of the propagator $[q_x(\chi_t|\chi_{t-1})]$ [18], we simply assume that there exists a unique stationary distribution $\pi_x$ such that $\pi(\chi_t) = \sum_{\chi_{t-1}} q_x(\chi_t|\chi_{t-1})\pi(\chi_{t-1})$. When given arbitrary initial distribution, evolution of memories finally goes to $\pi_x$ exponentially fast. Our conclusions are mainly drawn based on this stationary distribution, the *steady state*.

To characterize time irreversibility of $X$ in steady state, we define the steady state probability flux as the difference of the joint probabilities forward and backward in time,

$$J_x(\chi_t|\chi_{t-1}) = \pi_x(\chi_{t-1})q_x(\chi_t|\chi_{t-1}) - \pi_x(\widetilde{\chi}_t)q_x(\widetilde{\chi}_{t-1}|\widetilde{\chi}_t), \qquad (2)$$

where $\widetilde{\chi}$ denotes the time reversal of $\chi$, for example, $\widetilde{\chi}_t = [x_t, x_{t-1}, \ldots, x_{t-m+1}]$ (see Fig. 1). Here, "time-reverse memories" have to be emphasized because the time reversibility or time irreversibility of the system has memory dependence and the memories have to be flipped reversely when the process turns back in time. We can prove that $X$ is time irreversible if and only if $J_x = 0$ for all possible time-successive memories $\chi_{t-1}$ and $\chi_t$ (see Appendix A). The geometrical interpretation for $J_x$ that it is a tensor of order $(m+1)$ can be found in Appendix B.

It is noteworthy that for a Markovian system where the memory length $m = 1$, since $\chi_t = \widetilde{\chi}_t = x_t$ and $\chi_{t-1} = \widetilde{\chi}_{t-1} = x_{t-1}$, $J_x$ reduces to the well-known form $J_x(x_t|x_{t-1}) = \pi_x(x_{t-1})q_x(x_t|x_{t-1}) - \pi_x(x_t)q_x(x_{t-1}|x_t)$. The detailed balance equation $J_x = 0$ means that the Markovian system is in equilibrium steady state where all the system sequences can be flipped reversely along the time arrow with no costs. When $J_x = 0$ in a non-Markovian system, i.e., $\pi_x(\chi_{t-1})q_x(\chi_t|\chi_{t-1}) = \pi_x(\widetilde{\chi}_t)q_x(\widetilde{\chi}_{t-1}|\widetilde{\chi}_t)$, the system is time symmetrical between $\chi_{t-1}$ and the time-reverse memory $\widetilde{\chi}_t$, for the description of the time reversibility is equivalent to the equilibrium condition of a non-Markovian system, analogous to a Markov system. Consequentially, $\pi_x$ becomes the equilibrium stationary distribution of the memories in this equilibrium case.

In the framework of landscape-flux theory [11–13], the transition probability works as the driving force of a Markov system and the stationary distribution $\pi$ quantifies the landscape of the potential field in the probabilistic dynamics. The driving force (transition probability) $q$ can be decomposed into the landscape gradient and local flux velocity which correspond to the time-reversible and -irreversible parts of the dynamics, respectively. This quantitative picture fully depicts the underlying mechanics behind time-irreversible or nonequilibrium behavior of a Markov system. Here, we show that this quantitative characteristic also works when a non-Markovian system has (finite length) memory dependence.

We decompose the transition probabilities $q_x$ into two parts correspondingly as

$$q_x(\chi_t|\chi_{t-1}) = D_x(\chi_t|\chi_{t-1}) + B_x(\chi_t|\chi_{t-1}), \qquad (3)$$

with

$$D_x(\chi_t|\chi_{t-1}) = \frac{1}{2\pi_x(\chi_{t-1})}[\pi_x(\chi_{t-1})q_x(\chi_t|\chi_{t-1}) + \pi_x(\widetilde{\chi}_t)q_x(\widetilde{\chi}_{t-1}|\widetilde{\chi}_t)],$$

$$B_x(\chi_t|\chi_{t-1}) = \frac{1}{2\pi_x(\chi_{t-1})}J_x(\chi_t|\chi_{t-1}).$$

Here, $D_x$ satisfies $D_x(\chi_t|\chi_{t-1}) \geqslant 0$ and $\sum_{\chi_t} D_x(\chi_t|\chi_{t-1}) = 1$, i.e., $D_x(\chi_t|\chi_{t-1})$ can be regarded as a series of transition probabilities (or driving force) from memory $\chi_{t-1}$ to memory $\chi_t$. We can see that $\sum_{\chi_{t-1}} \pi_x(\chi_{t-1})D_x(\chi_t|\chi_{t-1}) = \pi_x(\chi_t)$, which means $\pi_x$ is the steady distribution corresponding to the transition probabilities $D_x$. Also, it can be verified that the detailed balance equation $\pi_x(\chi_{t-1})D_x(\chi_t|\chi_{t-1}) = \pi_x(\widetilde{\chi}_t)D_x(\widetilde{\chi}_{t-1}|\widetilde{\chi}_t)$ holds for $D_x$ and $\pi_x$. This means that the $X$ process behaves in a time-reversible way under the driving force $D_x$. Thus, $D_x$ can be expressed as $D_x(\chi_t|\chi_{t-1}) = D_x(\widetilde{\chi}_{t-1}|\widetilde{\chi}_t) \exp\{\log \pi_x(\widetilde{\chi}_t) - \log \pi_x(\chi_{t-1})\}$, where $\exp\{\log \pi_x(\widetilde{\chi}_t)\}$ can be viewed as the exponential of the difference between the potential $-\log \pi_x(\chi_{t-1})$ and $-\log \pi_x(\widetilde{\chi}_t)$. Then, $D_x$ can be identified as a gradientlike force corresponding to the landscape $\pi_x$ and the detailed balance (time-reversible) part of the stochastic information dynamics.

$B_x$ is the time-irreversible part (local flux velocity) of the stochastic information dynamics. This is because $B_x$ is not preserved time-reversal transformation. On the other hand, $B_x$ can be obtained directly from the probability flux $J_x$. A nonzero $B_x$ can be viewed as the force that drives the system to rotate among states with steady probability fluxes and to depart from detailed balance (or equilibrium state). The detailed balance breaking is indicated by $\pi_x(\chi_{t-1})B_x(\chi_t|\chi_{t-1}) = -\pi_x(\widetilde{\chi}_t)B_x(\widetilde{\chi}_{t-1}|\widetilde{\chi}_t)$.

### B. Interaction between two non-Markovian systems

When $X$ is interacting with another system $S$, we assume that $S$ generates a a finite-state, discrete-time, and irreducible chain with state space $\mathcal{S} = \{1, 2, \ldots, n\}$. We also assume that $S$ has $m$-dimensional memories $\varsigma_t = [s_{t-m+1}, s_{t-m+2}, \ldots, s_t]$ and transition probabilities $q_s(\varsigma_t|\varsigma_{t-1}) = q_s(s_t|s_{t-m}, s_{t-m+1}, \ldots, s_{t-1})$ for two successive memories $\varsigma_{t-1}$ and $\varsigma_t$. The evolution of distribution of its memories can be expressed by a Markov-type master equation

$$\Pr(\varsigma_t) = \sum_{\varsigma_{t-1}} q_s(\varsigma_t|\varsigma_{t-1})\Pr(\varsigma_{t-1}), \qquad (4)$$

with unique solution of stationary distribution $\pi_s$ which satisfies $\pi_s(\varsigma_t) = \sum_{\varsigma_{t-1}} q_s(\varsigma_t|\varsigma_{t-1})\pi_s(\varsigma_{t-1})$.

As we will show in Sec. V, the memories (and dynamics) of $X$ and $S$ are actually determined by the interactions which are fully embodied in the dynamics of the composite system $(X, S)$. Dynamics of $(X, S)$ can be obtained from the observations of the composite process. For the simplicity of discussion, we let $X$, $S$, and the composite system $(X, S)$ have

the same memory length. And this is also helpful to build up a unified framework for interacting non-Markovian systems with an explicit and computable structure. Here, we let the $m$-dimensional memories and transition probabilities of $(X, S)$ be $(\chi_t, \varsigma_t) = [(x_{t-m+1}, s_{t-m+1}), \ldots, (x_t, s_t)]$ and $q(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1}) = q(x_t, s_t|(x_{t-m}, s_{t-m}), \ldots, (s_{t-1}, s_{t-1}))$, respectively. Also, a Markov-type master equation for the composite system $(X, S)$ can be given by

$$\Pr(\chi_t, \varsigma_t) = \sum_{(\chi_{t-1}, \varsigma_{t-1})} q(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1})\Pr(\chi_{t-1}, \varsigma_{t-1}),$$

$$(5)$$

with unique solution of stationary distribution $\pi$.

For no loss of generality, we let $(X, S)$ be in its steady state with stationary distribution $\pi$ with both $X$ and $S$ achieving their own steady states with stationary distributions $\pi_x$ and $\pi_s$, respectively (jointly stationary assumption). This means the stationary distributions $\pi_x$ and $\pi_s$ corresponding to master equations (1) and (4) are the marginal distributions of $\pi$ corresponding to Eq. (5), i.e., both $\pi_s(\varsigma_t) = \sum_{\chi_t} \pi(\chi_t, \varsigma_t)$ and $\pi_x(\chi_t) = \sum_{\varsigma_t} \pi(\chi_t, \varsigma_t)$ hold in the steady state of the composite system. Similar to $J_x$ in Eq. (2), we define probability fluxes $J_s$ and $J$ for systems $S$ and $(X, S)$, respectively, for characterizing time irreversibility of corresponding processes. The landscape-flux decomposition can be also applied for transition probabilities of $S$ and $(X, S)$, respectively, in the form of Eq. (3): $q_s = D_s + B_s$ and $q = D + B$.

Due to jointly stationary assumption, relations between $J_x$, $J_s$, and $J$ can be shown as $J_x(\chi_t|\chi_{t-1}) = \sum_{\varsigma_{t-1}} \sum_{\varsigma_t} J(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1})$ and $J_s(\varsigma_t|\varsigma_{t-1}) = \sum_{\chi_{t-1}} \sum_{\chi_t} J(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1})$. These relations indicate that $J_x$ and $J_s$ can be regarded as two coarse-grained versions of $J$ with respect to subsystems $X$ and $S$, respectively. Also from these relations, we see that if $(X, S)$ is time reversible ($J = 0$), then both $X$ and $S$ are time reversible ($J_x = 0$ and $J_s = 0$). Conversely, if $X$ or $S$ is time irreversible ($J_x \neq 0$ or $J_s \neq 0$), then $(X, S)$ must be time irreversible ($J \neq 0$).

## III. NONEQUILIBRIUM THERMODYNAMICS OF INTERACTING NON-MARKOVIAN SYSTEMS

### A. Decomposition of mutual information rate into time-reversible and -irreversible parts and interactions between non-Markovian systems

As a fundamental concept in information theory, mutual information rate (MIR) [19] depicts the amount of information that two systems share. Here, we use MIR to quantify interaction between the non-Markovian systems $X$ and $S$. The MIR is defined by

$$I(X, S) = \lim_{T \to \infty} \frac{1}{T} \sum_{\Gamma_x(T), \Gamma_s(T)} \Pr(\Gamma_x(T), \Gamma_s(T))$$

$$\times \log \frac{\Pr(\Gamma_x(T), \Gamma_s(T))}{\Pr(\Gamma_x(T))\Pr(\Gamma_s(T))},$$

where time sequences $\Gamma_x(T) = [x_1, x_2, \ldots, x_T]$, $\Gamma_s(T) = [s_1, s_2, \ldots, s_T]$, and $(\Gamma_x(T), \Gamma_s(T))$ are generated by $X$, $S$, and $(X, S)$, respectively. The MIR is always non-negative and

it vanishes if and only if $X$ and $S$ do not interact with each other efficiently.

Due to the Markov-type master equations (1), (4), and (5), and given the initial distributions of memories to be stationary distributions, probabilities of $\Gamma_x(T)$, $\Gamma_s(T)$, and $(\Gamma_x(T), \Gamma_s(T))$ can be given by $\Pr(\Gamma_x(T)) = \pi_x(\chi_m) \prod_{i=m}^{T-1} q_x(\chi_{i+1}|\chi_i)$, $\Pr(\Gamma_s(T)) = \pi_s(\varsigma_m) \prod_{i=m}^{T-1} q_s(\varsigma_{i+1}|\varsigma_i)$, and $\Pr(\Gamma_x(T), \Gamma_s(T)) = \pi(\chi_m, \varsigma_m) \prod_{i=m}^{T-1} q(\chi_{i+1}, \varsigma_{i+1}|\chi_i, \varsigma_i)$, respectively.

We then have the explicit form of $I(X, S)$ which reads as (see the Appendix C)

$$I(X, S) = \sum_{(\chi_t, \varsigma_t)} \sum_{(\chi_{t-1}, \varsigma_{t-1})} \pi(\chi_{t-1}, \varsigma_{t-1}) q(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1})$$
$$\times\, i(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1}), \qquad (6)$$

where

$$i(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1}) = \log \frac{q(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1})}{q_x(\chi_t|\chi_{t-1}) q_s(\varsigma_t|\varsigma_{t-1})}.$$

Here, $i$ can be considered as the detailed interactions between $X$ and $S$ when a transition $(\chi_{t-1}, \varsigma_{t-1}) \rightarrow (\chi_t, \varsigma_t)$ occurs. And $I(X, S)$ is the average of $i$ all over the transitions of memories.

Corresponding to the landscape-flux decomposition of the driving forces, $I(X, S)$ can be decomposed into two parts:

$$I(X, S) = I_D(X, S) + I_B(X, S), \qquad (7)$$

where

$$I_D(X, S) = \sum_{(\chi_t, \varsigma_t)} \sum_{(\chi_{t-1}, \varsigma_{t-1})} \pi(\chi_{t-1}, \varsigma_{t-1}) D(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1})$$
$$\times\, i(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1}),$$

$$I_B(X, S) = \sum_{(\chi_t, \varsigma_t)} \sum_{(\chi_{t-1}, \varsigma_{t-1})} \pi(\chi_{t-1}, \varsigma_{t-1}) B(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1})$$
$$\times\, i(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1})$$
$$= \frac{1}{2} \sum_{(\chi_t, \varsigma_t)} \sum_{(\chi_{t-1}, \varsigma_{t-1})} J(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1})$$
$$\times\, i(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1})$$

are the time-reversible (detailed balance preserving) and -irreversible (detailed balance breaking) parts of the MIR, respectively. We give the physical interpretations to $I_D(X, S)$ and $I_B(X, S)$. We define the stochastic interactions between two possible time sequences of $X$ and $S$ in time $T$ as

$$k(\Gamma_x(T), \Gamma_s(T)) = \log \frac{\Pr(\Gamma_x(T), \Gamma_s(T))}{\Pr(\Gamma_x(T)) P(\Gamma_s(T))}$$
$$= \sum_{(\chi_t, \varsigma_t)} \sum_{(\chi_{t-1}, \varsigma_{t-1})} N(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1})$$
$$\times\, i(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1}),$$

where $N(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1})$ counts the number of transitions of $(\chi_{t-1}, \varsigma_{t-1}) \rightarrow (\chi_t, \varsigma_t)$ along the sequence $(\Gamma_x(T),$ $\Gamma_s(T))$, and $\lim_{T\to\infty} \frac{1}{T} \langle N(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1}) \rangle_{(\Gamma_x(T), \Gamma_s(T))} = \pi(\chi_{t-1}, \varsigma_{t-1}) q(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1})$. Clearly, the averaged interactions in sequence space of $X$ and $S$ and in long time limit is the MIR:

$\lim_{T\to\infty} \frac{1}{T} \langle k(\Gamma_x(T), \Gamma_s(T)) \rangle_{(\Gamma_x(T), \Gamma_s(T))} = I(X, S)$. We also define the stochastic interactions between two time reversals of the sequences of $X$ and $S$ as

$$k(\widetilde{\Gamma}_x(T), \widetilde{\Gamma}_s(T)) = \log \frac{\Pr(\widetilde{\Gamma}_x(T), \widetilde{\Gamma}_s(T))}{\Pr(\widetilde{\Gamma}_x(T)) \Pr(\widetilde{\Gamma}_s(T))}$$
$$= \sum_{(\widetilde{\chi}_t, \widetilde{\varsigma}_t)} \sum_{(\widetilde{\chi}_{t-1}, \widetilde{\varsigma}_{t-1})} N(\widetilde{\chi}_{t-1}, \widetilde{\varsigma}_{t-1}|\widetilde{\chi}_t, \widetilde{\varsigma}_t)$$
$$\times\, i(\widetilde{\chi}_{t-1}, \widetilde{\varsigma}_{t-1}|\widetilde{\chi}_t, \widetilde{\varsigma}_t),$$

where $N(\widetilde{\chi}_{t-1}, \widetilde{\varsigma}_{t-1}|\widetilde{\chi}_t, \widetilde{\varsigma}_t)$ counts the number of time-reversal transitions of $(\widetilde{\chi}_t, \widetilde{\varsigma}_t) \rightarrow (\widetilde{\chi}_{t-1}, \widetilde{\varsigma}_{t-1})$ along $(\widetilde{\Gamma}_x(T),$ $\widetilde{\Gamma}_s(T))$, and $\lim_{T\to\infty} \frac{1}{T} \langle N(\widetilde{\chi}_{t-1}, \widetilde{\varsigma}_{t-1}|\widetilde{\chi}_t, \widetilde{\varsigma}_t) \rangle_{(\Gamma_x(T), \Gamma_s(T))} = \lim_{T\to\infty} \frac{1}{T} \langle N(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1}) \rangle_{(\Gamma_x(T), \Gamma_s(T))}$. The change of stochastic interaction between the two subsystems when the time sequences turn back in time is quantified by the time-irreversible part of the stochastic interactions $k(\widetilde{\Gamma}_x(T), \widetilde{\Gamma}_s(T))$:

$$k_B(\Gamma_x(T), \Gamma_s(T)) = \tfrac{1}{2}(k(\Gamma_x(T), \Gamma_s(T)) - k(\widetilde{\Gamma}_x(T), \widetilde{\Gamma}_s(T))).$$

Clearly, $k_B$ measures the increasing ($k_B < 0$) or decreasing ($k_B > 0$) interactions between $X$ and $S$ along the time reversal $(\widetilde{\Gamma}_x(T), \widetilde{\Gamma}_s(T))$ compared to that of $(\Gamma_x(T), \Gamma_s(T))$. The time-reversible part of the $k(\widetilde{\Gamma}_x(T), \widetilde{\Gamma}_s(T))$ shows the remaining amount of interactions in both $(\Gamma_x(T), \Gamma_s(T))$ and $(\widetilde{\Gamma}_x(T), \widetilde{\Gamma}_s(T))$:

$$k_D(\Gamma_x(T), \Gamma_s(T)) = \tfrac{1}{2}(k(\Gamma_x(T), \Gamma_s(T)) + k(\widetilde{\Gamma}_x(T), \widetilde{\Gamma}_s(T))).$$

Thus, we have

$$I(X, S) = \lim_{T\to\infty} \frac{1}{T} \langle k(\Gamma_x(T), \Gamma_s(T)) \rangle_{\Gamma_z(T)}$$
$$= \lim_{T\to\infty} \frac{1}{T} \langle k_D(\Gamma_x(T), \Gamma_s(T)) \rangle_{\Gamma_z(T)}$$
$$+ \lim_{T\to\infty} \frac{1}{T} \langle k_B(\Gamma_x(T), \Gamma_s(T)) \rangle_{\Gamma_z(T)}$$
$$= I_D(X, S) + I_B(X, S),$$

where $I_D(X, S) = \lim_{T\to\infty} \frac{1}{T} \langle k_D(\Gamma_x(T), \Gamma_s(T)) \rangle_{\Gamma_z(T)}$, and $I_B(X, S) = \lim_{T\to\infty} \frac{1}{T} \langle k_B(\Gamma_x(T), \Gamma_s(T)) \rangle_{\Gamma_z(T)}$ with explicit forms being shown in Eq. (7).

Intuitively, if the two subsystems do not interact with each other efficiently [$I(X, S) = 0$], then we must have both $I_B(X, S) = 0$ and $I_D(X, S) = 0$. However, $I_B(X, S) = 0$ does not imply that $I(X, S) = 0$ because it related to not only the interactions, but also to time irreversibility of the composite systems and subsystems.

### B. Nonequilibrium thermodynamics of interacting non-Markovian systems: Relations between entropy production rates and mutual information rate

It is interesting to study the nonequilibrium thermodynamics of two interacting non-Markovian systems. To do so,

we explore the relation between the the mutual information rate and the entropy production rates of interacting systems. Entropy production rate (EPR) [20] of an information system characterizes the degree of the time-irreversible thermal information flows from the system to the nonequilibrium environments or, say, the rate of thermal information dissipation of the system under nonequilibrium conditions. It is defined by

$$R = \lim_{T \to \infty} \frac{1}{T} \sum_{\Gamma(T)} \Pr(\Gamma(T)) \log \frac{\Pr(\Gamma(T))}{\Pr(\widetilde{\Gamma}(T))} \geqslant 0,$$

for a stochastic system in steady state. The system is time reversible if and only if $R = 0$ or the system process is time irreversible, or we always have $R > 0$. Needless to say, the non-negativity of the EPR is not only a mathematical result, but also a key element of the second law of thermodynamics for open information systems: information (carried by energy and matter) is always dissipated irreversibly from the system with small degree of freedom to the environments with large degrees of freedom.

By noting master equations (1), (4), and (5), we realize that systems $X$, $S$, and $(X, S)$ should have Markov-type EPRs [21] with explicit form taking into account of the memory dependence although they may be non-Markovian systems. The corresponding EPRs can be given by (see Appendix C)

$$R(X) = \frac{1}{2} \sum_{\chi_t} \sum_{\chi_{t-1}} J_x(\chi_t | \chi_{t-1}) \log \frac{q_x(\chi_t | \chi_{t-1})}{q_x(\widetilde{\chi}_{t-1} | \widetilde{\chi}_t)},$$

$$R(S) = \frac{1}{2} \sum_{\varsigma_t} \sum_{\varsigma_{t-1}} J_s(\varsigma_t | \varsigma_{t-1}) \log \frac{q_s(\varsigma_t | \varsigma_{t-1})}{q_s(\widetilde{\varsigma}_{t-1} | \widetilde{\varsigma}_t)},$$

$$R(X, S) = \frac{1}{2} \sum_{(\chi_t, \varsigma_t)} \sum_{(\chi_{t-1}, \varsigma_{t-1})} J(\chi_t, \varsigma_t | \chi_{t-1}, \varsigma_{t-1})$$
$$\times \log \frac{q(\chi_t, \varsigma_t | \chi_{t-1}, \varsigma_{t-1})}{q(\widetilde{\chi}_{t-1}, \widetilde{\varsigma}_{t-1} | \widetilde{\chi}_t, \widetilde{\varsigma}_t)}. \quad (8)$$

Due to interactions between $X$ and $S$, the relation between $R(X)$, $R(S)$, and $R(X, S)$ can be revealed naturally by the time-irreversible part of MIR, $I_B(X, S)$, which gives

$$I_B(X, S)$$
$$= \frac{1}{2} \sum_{(\chi_t, \varsigma_t)} \sum_{(\chi_{t-1}, \varsigma_{t-1})} J(\chi_t, \varsigma_t | \chi_{t-1}, \varsigma_{t-1}) i(\chi_t, \varsigma_t | \chi_{t-1}, \varsigma_{t-1})$$
$$= \frac{1}{4} \sum_{(\chi_t, \varsigma_t)} \sum_{(\chi_{t-1}, \varsigma_{t-1})} J(\chi_t, \varsigma_t | \chi_{t-1}, \varsigma_{t-1}) \{ i(\chi_t, \varsigma_t | \chi_{t-1}, \varsigma_{t-1})$$
$$- i(\widetilde{\chi}_{t-1}, \widetilde{\varsigma}_{t-1} | \widetilde{\chi}_t, \widetilde{\varsigma}_t) \}$$
$$= \frac{1}{2} [R(X, S) - R(X) - R(S)]. \quad (9)$$

This equality provides another interpretation of the time-irreversible MIR. To see this, we focus on the differences

$$R(X|S) \equiv R(X, S) - R(S), \quad R(S|X) \equiv R(X, S) - R(X).$$
$$(10)$$

These differences can be viewed as the EPR (information dissipation) of one subsystem controlled by (or conditioning on) the time-irreversible behavior of the other subsystem. This is because, as shown in the definition of the EPR,

$$R(X|S) = R(X, S) - R(S)$$
$$= \lim_{T \to \infty} \frac{1}{T} \sum_{\Gamma_x(T), \Gamma_s(T)} \Pr(\Gamma_x(T), \Gamma_s(T))$$
$$\times \left\{ \log \frac{\Pr(\Gamma_x(T), \Gamma_s(T))}{\Pr(\widetilde{\Gamma_x}(T), \widetilde{\Gamma_s}(T))} - \log \frac{\Pr(\Gamma_s(T))}{\Pr(\widetilde{\Gamma_s}(T))} \right\}$$
$$= \lim_{T \to \infty} \frac{1}{T} \sum_{\Gamma_s(T)} \Pr(\Gamma_s(T)) \sum_{\Gamma_x(T)} \Pr(\Gamma_x(T) | \Gamma_s(T))$$
$$\times \log \frac{\Pr(\Gamma_x(T) | \Gamma_s(T))}{\Pr(\widetilde{\Gamma_x}(T) | \widetilde{\Gamma_s}(T))},$$

where $(\Gamma_x, \Gamma_s)$ $\Gamma_x$ and $\Gamma_s$ denote the possible time sequences of the composite system and subsystems, respectively; $\widetilde{\Gamma}$ denotes the corresponding time-reverse time sequence. Here, the conditional probabilities $\Pr(\Gamma_x(T) | \Gamma_s(T)) = \frac{\Pr(\Gamma_x(T), \Gamma_s(T))}{\Pr(\Gamma_s(T))}$ and $\Pr(\widetilde{\Gamma_x}(T) | \widetilde{\Gamma_s}(T)) = \frac{\Pr(\widetilde{\Gamma_x}(T), \widetilde{\Gamma_s}(T))}{\Pr(\widetilde{\Gamma_s}(T))}$ reveal the time-forward and -backward behaviors of $X$ which are controlled by the the time-forward and -backward behaviors of $S$ correspondingly. Similarly, we can derive the expression of $R(S|X)$ as follows:

$$R(S|X) = R(X, S) - R(X)$$
$$= \lim_{T \to \infty} \frac{1}{T} \sum_{\Gamma_x(T)} \Pr(\Gamma_x(T)) \sum_{\Gamma_s(T)} \Pr(\Gamma_s(T) | \Gamma_x(T))$$
$$\times \log \frac{\Pr(\Gamma_s(T) | \Gamma_x(T))}{\Pr(\widetilde{\Gamma_s}(T) | \widetilde{\Gamma_x}(T))}.$$

Then, the equality with respect to the time-irreversible MIR $I_B(X, S)$ and the EPRs in Eq. (9) can be rearranged as follows:

$$R(X|S) = 2I_B(X, S) + R(X), \quad R(S|X) = 2I_B(X, S) + R(S).$$
$$(11)$$

Thus, the information dissipation of one subsystem [$R(X|S)$ or $R(S|X)$] controlled by the time-irreversible behavior of the other subsystem (or a time-variant environment) is constituted by the "self"-information dissipation of the subsystem [$R(X)$ or $R(S)$] and the information dissipation associating with the interaction $I_B(X, S)$. Here, "information dissipation associating with the interaction" does not mean that $I_B(X, S)$ has to be non-negative. In fact, if $R(X|S) > R(X)$ or $R(S|X) > R(S)$, then we obtain a positive $I_B(X, S)$ which means the information of interaction is dissipated from the systems into the environments time irreversibly. Otherwise, a negative $I_B(X, S)$ [$R(X|S) < R(X)$ or $R(S|X) < R(S)$] means that the information of interaction is dissipated from the the environments into the systems to maintain the self-information dissipation of the systems. This fully provides an understanding of the intrinsic property of the time-irreversible MIR. A related discussion about time-reversed control can be found in [22].

## IV. NON-MARKOV SYSTEMS WITH UNKNOWN MEMORY LENGTHS

### A. Finite memory approximation

Quite often, we do not always know the memory lengths of most non-Markovian systems even though they are elaborately constructed. Besides, there are two disadvantages on analyzing systems with unknown memory lengths: (1) the two interacting systems may have different memory lengths; (2) any one of the subsystems may have a really large memory length. However, by using the method of finite memory approximation (FMA), we may easily include more general class of non-Markovian cases into the framework of landscape-flux theory for information dynamics.

A non-Markovian system often emerges when it is interacting with another system with comparable size and state-switching rate. While the interaction is identified, memories of both interacting systems are determined. This interaction can be fully embodied in the dynamics of composite systems. This means we can obtain complete dynamics of non-Markov subsystems from composite system. Thus, we consider a composite system $Z = (X, S)$ with Markov-type master equation shown in Eq. (5) and the memory lengths of the subsystems being unknown. Here, we use $m_z$ to denote the exact memory length of $Z = (X, S)$ in Eq. (5). From FMA, by choosing a unified memory length for the two subsystems we can approximate the processes of the subsystems in both stochastic dynamics and thermodynamics with arbitrary precision. This is the direct result of the martingale convergence theorem [17].

We assume that the subsystems $X$ and $S$ satisfy the jointly stationary assumption. We then let $\mathfrak{X}^{(m)}$ and $\mathfrak{S}^{(m)}$ be the approximating systems of $X$ and $S$, respectively. Here, $m \geqslant m_z$ is the unified memory length of both $\mathfrak{X}$ and $\mathfrak{S}$. Then, the transition probabilities of successive memories of $\mathfrak{X}^{(m)}$ and $\mathfrak{S}^{(m)}$ can be evaluated by

$$q_x\big(\chi_t^{(m)}\big|\chi_{t-1}^{(m)}\big) = \frac{\sum_{\varsigma_t^{(m)}} \Pr\big(\chi_t^{(m)}, \varsigma_t^{(m)}\big)}{\sum_{\varsigma_t^{(m+1)}} \Pr\big(\chi_t^{(m+1)}, \varsigma_t^{(m+1)}\big)},$$

$$q_s\big(\varsigma_t^{(m)}\big|\varsigma_{t-1}^{(m)}\big) = \frac{\sum_{\chi_t^{(m)}} \Pr\big(\chi_t^{(m)}, \varsigma_t^{(m)}\big)}{\sum_{\chi_t^{(m+1)}} \Pr\big(\chi_t^{(m+1)}, \varsigma_t^{(m+1)}\big)},$$

where $\chi_t^{(m)} = [x_{t-m+1}, \ldots, x_t]$ is a sequence with superscript $(m)$ denoting its length and subscript denoting its end time. Similar settings are given to the sequences (or memories) appearing in this section. The joint probabilities $\Pr(\chi_t^{(m)}, \varsigma_t^{(m)})$ and $\Pr(\chi_t^{(m+1)}, \varsigma_t^{(m+1)})$ are given by dynamics of $Z$ in steady state:

$$\Pr\big(\chi_t^{(m)}, \varsigma_t^{(m)}\big) = \pi\big(\chi_{t-m+m_z}^{(m_z)}, \varsigma_{t-m+m_z}^{(m_z)}\big) \prod_{i=t-m+m_z+1}^{t}$$
$$\times q\big(\chi_i^{(m_z)}, \varsigma_i^{(m_z)}\big|\chi_{i-1}^{(m_z)}, \varsigma_{i-1}^{(m_z)}\big),$$

$$\Pr\big(\chi_t^{(m+1)}, \varsigma_t^{(m+1)}\big) = \pi\big(\chi_{t-m+m_z-1}^{(m_z)}, \varsigma_{t-m+m_z-1}^{(m_z)}\big) \prod_{i=t-m+m_z}^{t}$$
$$\times q\big(\chi_i^{(m_z)}, \varsigma_i^{(m_z)}\big|\chi_{i-1}^{(m_z)}, \varsigma_{i-1}^{(m_z)}\big).$$

Then, we can formulate Markov-type master equations for $\mathfrak{X}^{(m)}$ and $\mathfrak{S}^{(m)}$, respectively, and have

$$\Pr\big(\chi_t^{(m)}\big) = \sum_{\chi_{t-1}^{(m)}} q_x\big(\chi_t^{(m)}\big|\chi_{t-1}^{(m)}\big)\Pr\big(\chi_{t-1}^{(m)}\big),$$

$$\Pr\big(\varsigma_t^{(m)}\big) = \sum_{\varsigma_{t-1}^{(m)}} q_s\big(\varsigma_t^{(m)}\big|\varsigma_{t-1}^{(m)}\big)\Pr\big(\varsigma_{t-1}^{(m)}\big).$$

Thus, $\mathfrak{X}^{(m)}$ and $\mathfrak{S}^{(m)}$ form stationary and ergodic processes when initial distributions are $\pi_x(\chi_t^{(m)}) = \sum_{\varsigma_t^{(m)}} \Pr(\chi_t^{(m)}, \varsigma_t^{(m)})$ and $\pi_s(\varsigma_t^{(m)}) = \sum_{\chi_t^{(m)}} \Pr(\chi_t^{(m)}, \varsigma_t^{(m)})$, respectively. For calculations of FMA, we need to increase the memory length of $Z$ from $m_z$ to $m$, and construct a pseudodynamics of $Z$ in steady state:

$$\Pr\big(\chi_t^{(m)}, \varsigma_t^{(m)}\big) = \sum_{(\chi_{t-1}^{(m)}, \varsigma_{t-1}^{(m)})} q\big(\chi_t^{(m)}, \varsigma_t^{(m)}\big|\chi_{t-1}^{(m)}, \varsigma_{t-1}^{(m)}\big)$$
$$\times \Pr\big(\chi_{t-1}^{(m)}, \varsigma_{t-1}^{(m)}\big),$$

where $q(\chi_t^{(m)}, \varsigma_t^{(m)}|\chi_{t-1}^{(m)}, \varsigma_{t-1}^{(m)}) = q(\chi_t^{(m_z)}, \varsigma_t^{(m_z)}|\chi_{t-1}^{(m_z)}, \varsigma_{t-1}^{(m_z)})$.

Now, we have a complete construction of dynamics with memory length $m$ similar to Eqs. (1), (4), and (5). We then evaluate the probability fluxes $J_x$, $J_s$, and $J$ for the approximating systems according to Eq. (2). Landscape-flux decomposition can be made for the driving forces of the approximating systems according to Eq. (3) with the steady state distributions quantifying the information landscapes while the steady state probability fluxes measuring the degree of the nonequilibriumness. In thermodynamics, entities become functions of $m$, they are denoted by $I(\mathfrak{X}^{(m)}, \mathfrak{S}^{(m)})$ (the MIR), $I_B(\mathfrak{X}^{(m)}, \mathfrak{S}^{(m)})$ (the time-irreversible part of MIR), $R(\mathfrak{X}^{(m)})$ (the EPR of $X$), and $R(\mathfrak{S}^{(m)})$ (the EPR of $S$), etc. These entities can be evaluated by using Eqs. (6), (7), and (8), respectively.

If $X$ and $S$ have different but finite memory lengths $m_x$ and $m_s$, respectively, we let the unified memory length $m \geqslant \max\{m_x, m_s, m_z\}$ and have

$$q_x\big(\chi_t^{(m)}\big|\chi_{t-1}^{(m)}\big) = q_x(x_t|x_{t-m}, \ldots, x_{t-1}) = q_x(x_t|x_{t-m_x}, \ldots, x_{t-1}) = q_x\big(\chi_t^{(m_x)}\big|\chi_{t-1}^{(m_x)}\big),$$

$$q_s\big(\varsigma_t^{(m)}\big|\varsigma_{t-1}^{(m)}\big) = q_s(s_t|s_{t-m}, \ldots, s_{t-1}) = q_s(s_t|s_{t-m_s}, \ldots, s_{t-1}) = q_s\big(\varsigma_t^{(m_s)}\big|\varsigma_{t-1}^{(m_s)}\big),$$

$$q\big(\chi_t^{(m)}, \varsigma_t^{(m)}\big|\chi_{t-1}^{(m)}, \varsigma_{t-1}^{(m)}\big) = q\big(\chi_t^{(m_z)}, \varsigma_t^{(m_z)}\big|\chi_{t-1}^{(m_z)}, \varsigma_{t-1}^{(m_z)}\big),$$

$$\sum_{x_{t-m},\ldots,x_{t-m_x-1}} J_x\big(\chi_t^{(m)}\big|\chi_{t-1}^{(m)}\big) = J_x\big(\chi_t^{(m_x)}\big|\chi_{t-1}^{(m_x)}\big),$$

$$\sum_{s_{t-m},\ldots,s_{t-m_s-1}} J_s\big(\varsigma_t^{(m)}\big|\varsigma_{t-1}^{(m)}\big) = J_s\big(\varsigma_t^{(m_s)}\big|\varsigma_{t-1}^{(m_s)}\big),$$

$$\sum_{(x_{t-m},s_{t-m}),\ldots,(x_{t-m_x-1},s_{t-m_s-1})} J\big(\chi_t^{(m)},\varsigma_t^{(m)}\big|\chi_{t-1}^{(m)},\varsigma_{t-1}^{(m)}\big) = J\big(\chi_t^{(m_z)},\varsigma_t^{(m_z)}\big|\chi_{t-1}^{(m_z)},\varsigma_{t-1}^{(m_z)}\big).$$

Hence, for $m \geqslant \max\{m_x, m_s, m_z\}$ we have

$$I(\mathfrak{X}^{(m)}, \mathfrak{S}^{(m)}) = I(X, S),$$

$$R(\mathfrak{X}^{(m)}) = R(X),$$

$$R(\mathfrak{S}^{(m)}) = R(S),$$

$$I_B(\mathfrak{X}^{(m)}, \mathfrak{S}^{(m)}) = I_B(X, S).$$

If $X$ and $S$ have unknown and large (and also different) memory lengths, according to the martingale convergence theorem as the memory length $m \to \infty$ the conditional entities

$$q_x(x_t|x_{t-m}, \ldots, x_{t-1}) \to q_x(x_t|x_{-\infty}, \ldots, x_{t-1}),$$

$$q_s(s_t|s_{t-m}, \ldots, s_{t-1}) \to q_s(s_t|s_{-\infty}, \ldots, s_{t-1}),$$

$$I(\mathfrak{X}^{(m)}, \mathfrak{S}^{(m)}) \to I(X, S),$$

$$R(\mathfrak{X}^{(m)}) \to R(X),$$

$$R(\mathfrak{S}^{(m)}) \to R(S),$$

$$I_B(\mathfrak{X}^{(m)}, \mathfrak{S}^{(m)}) \to I_B(X, S),$$

almost surely. This means we can use the FMA method to evaluate these entities with arbitrary precision.

For numerical calculations, for example, to evaluate $R(X)$ and $I_B(X, S)$, we can use the following sequences:

$$R(\mathfrak{X}^{(m_z)}), R(\mathfrak{X}^{(m_z+1)}), R(\mathfrak{X}^{(m_z+2)}), \ldots$$

and

$$I_B(\mathfrak{X}^{(m_z)}, \mathfrak{S}^{(m_z)}), I_B(\mathfrak{X}^{(m_z+1)}, \mathfrak{S}^{(m_z+1)}),$$
$$I_B(\mathfrak{X}^{(m_z+2)}, \mathfrak{S}^{(m_z+2)}) \ldots.$$

The numerical calculations will continue until

$$\left| \frac{R(\mathfrak{X}^{(M_1+1)}) - R(\mathfrak{X}^{(M_1)})}{R(\mathfrak{X}^{(M_1)})} \right| \leqslant \delta_1,$$

$$\left| \frac{I_B(\mathfrak{X}^{(M_2+1)}, \mathfrak{S}^{(M_2+1)}) - I_B(\mathfrak{X}^{(M_2)}, \mathfrak{S}^{(M_2)})}{I_B(\mathfrak{X}^{(M_2)}, \mathfrak{S}^{(M_2)})} \right| \leqslant \delta_2,$$

for given small enough thresholds of the relative errors $\delta_1$ and $\delta_2$ at $M_1$ and $M_2$, respectively. Then, we take $R(\mathfrak{X}^{(M_1)})$ and $I_B(\mathfrak{X}^{(M_2)}, \mathfrak{S}^{(M_2)})$ as the true values of $R(X)$ and $I_B(X, S)$, respectively. Here, $M_1$ needs not to be equal to $M_2$ because $M_1$ approximates to the true value of $m_x$ and $M_2$ measures the true value of $m = \max\{m_x, m_s, m_z\}$. Similar algorithm can be designed for other entities.

## V. BOUNDS OF NONEQUILIBRIUM ENTITIES

Although the FMA method provides the chance that we can obtain observables of a non-Markovian system with arbitrary precision, we should note that we may not obtain adequate information of the stochastic dynamics of the composite systems or subsystems in finite time especially when only the dynamics of the composite system could be observed. On the other hand, when the state spaces or the memory lengths are quite large, the complexity of calculations of the FMA method increases intolerably. It then becomes unworthy (or even impossible) to evaluate the exact values of the system observables under concern. In these situations, it would be meaningful to check whether the designed systems have the desired properties of entities within certain bounds rather than to evaluate the exact values for practical applications. Here, the desired properties refer to the characterization of system information dissipation (the EPRs) and characterization of interaction information dissipation (the time-irreversible MIR). These bounds involve the log-sum inequality which is useful to derive the bounds of entities related to the Kulback-Leibler divergence (see Appendix D).

### A. Upper bounds of entropy production rates of subsystems and lower bound of time-irreversible mutual information rate

For the EPRs of the subsystems, $R(X)$ and $R(S)$ in Eq. (8), the log-sum inequality suggests that

$$R(X, S) \geqslant R(X), \quad R(X, S) \geqslant R(S), \tag{12}$$

where $R(X, S) = R(X)$ holds for that $G = \frac{\pi(\chi_{t-1}, \varsigma_{t-1})q(\chi_t, \varsigma_t|\chi_{t-1}, \varsigma_{t-1})}{\pi(\tilde{\chi}_t, \tilde{\varsigma}_t)q(\tilde{\chi}_{t-1}, \tilde{\varsigma}_{t-1}|\tilde{\chi}_t, \tilde{\sigma}_t)}$ is not a function of $\varsigma_t$ and $\varsigma_{t-1}$, and $R(X, S) = R(S)$ holds for that $G$ is not a function of $\chi_t$ and $\chi_{t-1}$. Here, $\chi$ and $\varsigma$ refer to the memories of $X$ and $S$, respectively. $\pi$ and $q$ refer to the stationary distribution and transition probability of the composite system, respectively. The equalities mean that the information dissipation of the composite system only depends on the information dissipation of one subsystem. For example, if $X$ is independent of $S$ and $R(S) = 0$, then we have $R(X, S) = R(X)$.

The inequalities of EPRs follow the intuition that the information dissipation of any subsystems cannot exceed that of the composite system. They also imply that the EPRs under control, $R(X|S)$ and $R(S|X)$, are always non-negative:

$$R(X|S) = R(Z) - R(S) \geqslant 0, \ R(S|X) = R(Z) - R(X) \geqslant 0. \tag{13}$$

By combining these inequalities with Eq. (9), we have a lower bound of time-irreversible MIR $I_B(X, S)$, which reads as

$$I_B(X, S) \geqslant \tfrac{1}{2} \max\{-R(X), -R(S)\}. \qquad (14)$$

This provides a vital constraint on the information dissipation of the interaction. This inequality makes sure that the information dissipation of any subsystems cannot exceed the total information dissipation $R(X, S)$.

### B. Lower bounds of entropy production rates and upper bound of time-irreversible mutual information rate

The FMA allows us to obtain series of the lower bounds of EPRs and upper bounds of the time-irreversible MIR easily from the approximating systems. If, for example, we use the FMA to approximate a non-Markovian system $X$ in steady state via two memory lengths $m_1$ and $m_2$ ($m_1 \geqslant m_2$) and the corresponding approximating systems are denoted by $\mathfrak{X}^{(m_1)}$ and $\mathfrak{X}^{(m_2)}$, then the log-sum inequality suggests an inequality relation between the EPRs of the two approximations $R(\mathfrak{X}^{(m_1)}) \geqslant R(\mathfrak{X}^{(m_2)})$ (see Appendix D). Generally, this relation can be extended to

$$R(X) \geqslant \cdots \geqslant R(\mathfrak{X}^{(m_1)}) \geqslant R(\mathfrak{X}^{(m_2)}) \geqslant \cdots \geqslant R(\mathfrak{X}^{(1)}) \geqslant 0, \qquad (15)$$

where $1 \leqslant \cdots \leqslant m_2 \leqslant m_1 \leqslant \cdots \leqslant m$ with $m$ being the exact memory length of $X$. Here, $R(\mathfrak{X}^{(1)})$ refers to the Markovian approximation of $X$. This inequality provides the series of lower bounds of the EPR of a non-Markovian system. These lower bounds indicate that decreasing memory length causes decreasing rate of information dissipation. In other words, the memories with certain length can be viewed as the "time" environments of a non-Markovian information system. Larger memory length implies an environment with larger degree of freedom and larger information dissipation.

Although 0 is the natural lower bound of the EPR, we should note that a nonzero lower bound $R(\mathfrak{X}^{(M)}) > 0$ ($M \leqslant m$) can be more informative than the trivial 0. To say the least, if we are only interested in the time irreversibility of $X$, then we can choose a smaller memory length, such as $M = 1 < m$, and confirm that $R(\mathfrak{X}^{(M)}) > 0$. Then, $X$ is time irreversible.

By noting Eqs. (9) and (10), provided the exact memory length of the composite system $m_z$ and the corresponding EPR $R(Z)$, we have the upper bounds of $I_B(X, S)$ as

$$
\begin{aligned}
I_B(X, S) &\leqslant \cdots \leqslant \tfrac{1}{2}[R(Z) - R(\mathfrak{X}^{(m_2)}) - R(\mathfrak{S}^{(m_2)})] \\
&\leqslant \tfrac{1}{2}[R(Z) - R(\mathfrak{X}^{(m_1)}) - R(\mathfrak{S}^{(m_1)})] \leqslant \cdots \\
&\leqslant \tfrac{1}{2}[R(Z) - R(\mathfrak{X}^{(1)}) - R(\mathfrak{S}^{(1)})] \\
&\leqslant \tfrac{1}{2}R(Z). \qquad (16)
\end{aligned}
$$

These inequalities imply that increasing memory lengths of subsystems (with the EPR of the composite system being fixed) decreases the time-irreversible MIR $I_B$. Since self-information dissipation of subsystems decreases with decreasing memory, the dissipation of the interactions (time-irreversible MIR) increases. Moreover, smaller memory length (for instance, the Markovian approximation) can be quite effective to bound the time-irreversible MIR.
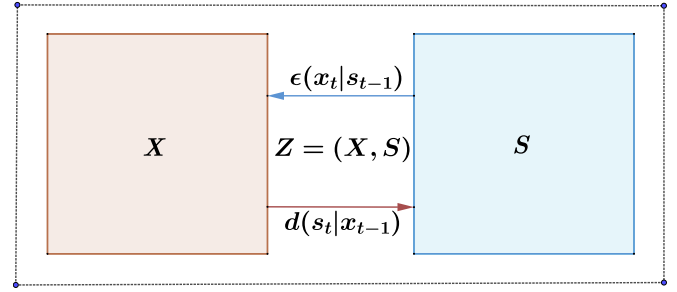


FIG. 2. Diagram of the two interacting information systems $X$ and $S$ in the example.

We should note that inequalities in Eqs. (12), (13), and (14) provide the general constraints on the EPRs of the subsystems and time-irreversible MIR.

## VI. AN EXAMPLE OF TWO INTERACTING INFORMATION SYSTEMS

To clarify the meaning of the information-theoretical formulations in the above, we construct a comprehensive example of two interacting information systems, which involves the information dissipation and (feedback) control. Related examples can be found in [10,23,24,25].

### A. Description of stochastic dynamics, analytical solutions of transition probabilities by using finite memory length method

In this example, two finite-state information systems denoted by $X$ and $S$, respectively, are interacting with each other. Their states are labeled by $\mathcal{X} = \{x = 1, 2, \ldots, l\}$ and $\mathcal{S} = \{s = 1, 2, \ldots, n\}$, respectively. Both systems are driven by several random information environments individually so that their behaviors become randomized. The interaction between them is considered to be purely informational without physical contact. Their stochastic dynamics can be described as follows (see Fig. 2).

The state of any one of the systems at $t$ is controlled by the state of the other system at $t - 1$. Since the systems are driven by random environments, then the controls become randomized which are characterized by the conditional probabilities $\epsilon(x_t|s_{t-1})$ and $d(s_t|x_{t-1})$. Here, $\epsilon$ indicates the probability of state of $X$ ($x_t$) controlled by state of $S$ ($s_{t-1}$). $d$ represents the probability of state of $S$ ($s_t$) controlled by $X$ ($x_{t-1}$). These conditional probabilities satisfy $\{\epsilon(x|s) : \epsilon(x|s) \geqslant 0, \sum_{x \in \mathcal{X}} \epsilon(x|s) = 1\}$ and $\{d(s|x) : d(s|x) \geqslant 0, \sum_{s \in \mathcal{S}} d(s|x) = 1\}$. The controls from the two systems are considered to be independent of each other at $t$. This means for the composite system $Z = (X, S)$, the transition probability from state at $t - 1$ to state at $t$ satisfies that $q(x_t, s_t|x_{t-1}, s_{t-1}) = \epsilon(x_t|s_{t-1})d(s_t|x_{t-1})$. Since the state of $Z$ at $t$ only depends on the state at $t - 1$, then $Z$ follows a Markovian dynamics.

Due to the description of the stochastic dynamics of $S$ and $X$, the probabilistic master equation with respect to $Z$ arises from the transition probabilities which reads as

$$\Pr(x_t, s_t) = \sum_{x_{t-1}, s_{t-1}} q(x_t, s_t|x_{t-1}, s_{t-1})\Pr(x_{t-1}, s_{t-1}).$$

We assume that there is a unique stationary distribution of $z$ denoted by $\pi$. We also assume that the $X$ and $S$ satisfy the joint stationary assumption. It can be seen from the description of the dynamics of $X$ and $S$ that the system states $x_t$ and $s_t$ are independent of each other. Thus, it is clear that the stationary marginal probabilities $\pi_x(x) = \sum_s \pi(x,s)$ and $\pi_s(s) = \sum_x \pi(x,s)$ satisfy $\pi(x,s) = \pi_x(x)\pi_s(s)$ (a proof of this can be found in Appendix E). This indicates that the time-sliced interaction (mutual information) $I(\mathfrak{X}, \mathfrak{S}) = \langle \log \frac{\pi(x,s)}{\pi_x(x)\pi_s(s)} \rangle = 0$. However, this does not mean the systems are independent of each other since $x_t$ and $s_t$ are controlled by $s_{-1}$ and $x_{-1}$, respectively.

In spite of specific conditional probabilities $\epsilon$ and $d$ for Markovian cases, $X$ and $S$ are both non-Markovian with unknown memory lengths in general. However, by applying the FMA method analytically, we have, when the approximating memory length $m \geqslant 2$ (see Appendix E),

$$q_x\big(\chi_t^{(m)}\big|\chi_{t-1}^{(m)}\big) = q_x\big(\chi_t^{(2)}\big|\chi_{t-1}^{(2)}\big) = \sum_{s_{t-1}} d(s_{t-1}|x_{t-2})\epsilon(x_t|s_{t-1})$$

$$= q_x(x_t|x_{t-2}),$$

$$q_s\big(\varsigma_t^{(m)}\big|\varsigma_{t-1}^{(m)}\big) = q_s\big(\varsigma_t^{(2)}\big|\varsigma_{t-1}^{(2)}\big) = \sum_{x_{t-1}} \epsilon(x_{t-1}|s_{t-2})d(s_t|x_{t-1})$$

$$= q_s(s_t|s_{t-2}).$$

These transition probabilities indicate that both $X$ and $S$ have memory length of 2. However, they can still be regarded as "sampled Markovian" systems since $q_x$ and $q_s$ have no correlation with the states at $t-1$. For example, $x_{t-2}$ exerts influence on $x_t$ via $s_{t-1}$ but not via $x_{t-1}$ directly. Thus, they can be rewritten as $q_x(x_t|x_{t-2})$ and $q_s(s_t|s_{t-2})$, respectively. But, neither system can be taken as Markovian systems because the knowledge of states at $t-2$ and $t-1$ must be known for generating the processes of both systems. The stationary distribution of $X$ and $S$ is recognized as $\pi_x(x)$ and $\pi_s(s)$ because the equalities $\pi_x(x_t) = \sum_{x_{t-2}} q_x(x_t|x_{t-2})\pi_x(x_{t-2})$ and $\pi_s(s_t) = \sum_{s_{t-2}} q_s(s_t|s_{t-2})\pi_s(s_{t-2})$ both hold in steady state.

It is noteworthy that the transition probabilities $q_x$ and $q_s$ can be written into the sums of the series of transition probabilities via different "channels," respectively,

$$q_x(x_t|x_{t-2}) = \sum_s q_x(x_t, s|x_{t-2}), \quad q_s(s_t|s_{t-2})$$

$$= \sum_x q_s(s_t, x|s_{t-2}),$$

where $q_x(x_t, s|x_{t-2}) = d(s|x_{t-2})\epsilon(x_t|s)$ ($s = s_{t-1}$ for short); $q_s(s_t, x|s_{t-2}) = \epsilon(x|s_{t-2})d(s_t|x)$ ($x = x_{t-1}$ for short); the notation forms of $q_x(x_t, s|x_{t-2})$ and $q_s(s_t, x|s_{t-2})$ are given by the chain rule of probabilities. By the description of dynamics of the systems, $s$ is determined by $x_{t-2}$ via $d(s|x_{t-2})$ then $x_t$ is determined by $s$ via $\epsilon(x_{t-2}|s)$. Thus, $q_x(x_t, s|x_{t-2})$ is the transition probability that $X$ jumps from $x_{t-2}$ to $x_t$ through the channel $s$. Analogously, $q_s(s_t, x|s_{t-2})$ is the transition probability that $S$ jumps from $s_{t-2}$ to $s_t$ through the channel $x$. These two series of transition probabilities would be helpful for clarifying the meaning of the time-irreversible MIR in this model.

TABLE I. An example of encoding the states $x_{t-1}$ and $x_t$ ($\mathcal{X} = \{1, 2, 3, 4\}$) by the Huffman coding. $x_{t-1}$ and $x_t$ are controlled by $s_{t-1} = 1$ and $s_{t-2} = 2$, respectively. The optimal code lengths $l_x$ are also shown in the table.

|  | $x_{t-1} = 1$ | $x_{t-1} = 2$ | $x_{t-1} = 3$ | $x_{t-1} = 4$ |
|---|---|---|---|---|
| $\epsilon(x_{t-1}|s_{t-2})$ | 1/4 | 1/4 | 1/4 | 1/4 |
| $l_x(x_{t-1}|s_{t-2})$ | 2 | 2 | 2 | 2 |
| Codewords | 01 | 00 | 11 | 10 |
|  | $x_t = 1$ | $x_t = 2$ | $x_t = 3$ | $x_t = 4$ |
| $\epsilon(x_t|s_{t-1})$ | 1/2 | 1/4 | 1/8 | 1/8 |
| $l_x(x_t|s_{t-1})$ | 1 | 2 | 3 | 3 |
| Codewords | 0 | 10 | 111 | 110 |

In this model, $-\log \pi_x(x)$ and $-\log \pi_s(s)$ can be taken as the potential landscapes of the subsystems $X$ and $S$, respectively. These two landscapes measure the self-information of the system states in bits which is independent of the time and the other system. The information fluxes of $X$ and $S$ can be taken as

$$J_x(x_t|x_{t-2}) = \pi_x(x_{t-2})q_x(x_t|x_{t-2}) - \pi_x(x_t)q_x(x_{t-2}|x_t),$$

$$J_s(s_t|s_{t-2}) = \pi_s(s_{t-2})q_s(s_t|s_{t-2}) - \pi_s(s_t)q_s(s_{t-2}|s_t).$$

### B. Dynamical observables of stochastic dynamics

Since the stochastic dynamics of the systems in this case is identified in the above, we can evaluate the nonequilibrium entities: the EPRs $R(Z)$, $R(X)$, and $R(S)$ by using Eq. (8), and the time-irreversible MIR $I_B(X, S)$ by using Eq. (7) analytically. In this model, we can give the nonequilibrium entities an intuitive picture by connecting them to the dynamical observables. According to the argument by Shannon [14], the system states can be encoded into series of codewords (sequences of 0 and 1) by using the optimal coding methods such as Shannon-Fano coding [15] or Huffman coding [26]. The (optimal) code lengths are equal to the negative logarithmic (conditional) probabilities corresponding to the system dynamics. These lengths measure the bits of information that the events with respect to system states (such as transitions) are observed by outer observers. These lengths can be observed stochastically and can be summed up or averaged along with the time sequences. Thus, they can be taken as the dynamical observables of systems. These dynamical observables (optimal lengths of codewords) can be classified into two groups: one group associating with the controls between systems and one group associating with the transitions between system states.

(1) The conditional probabilities $d$ and $\epsilon$ contain the detailed information of the (feedback) controls between systems. According to the argument by Shannon, we can assign a codeword for every state of $X$ with respect to a fixed control condition $s$ and assign a codeword for every state of $S$ with respect to a fixed control condition $x$. An example can be found in Table I. The optimal code lengths $l_x$ and $l_s$ are equal to the bits of information of the (feedback) controls between systems (negative logarithmic control probabilities)

correspondingly, which read as

$$l_x(x_t|s_{t-1}) \equiv -\log \epsilon(x_t|s_{t-1}),$$
$$l_s(s_t|x_{t-1}) \equiv -\log d(s_t|x_{t-1}),$$

where we use the convention $\log = \log_2$. The optimal code lengths quantify the detailed control information with bits. They can be taken as the dynamical observables associating with the controls.

(2) The transition probabilities with respect to the composite systems and subsystems $q$, $q_x$, and $q_s$ contain the detailed information of the detailed system dynamics. We can assign a codeword for every state transition corresponding to the transition probabilities. The optimal code lengths are equal to the bits of information of the transitions between system states (negative logarithmic transition probabilities) correspondingly, which read as

$$h_z(x_t, s_t|x_{t-1}, s_{t-1}) \equiv -\log q(x_t, s_t|x_{t-1}, s_{t-1}),$$
$$h_x(x_t|x_{t-2}) \equiv -\log q_x(x_t|x_{t-2}),$$
$$h_s(s_t|s_{t-2}) \equiv -\log q_s(s_t|s_{t-2}),$$
$$h_x^{(s)}(x_t|x_{t-2}) \equiv -\log q_x(x_t, s|x_{t-2}),$$
$$h_s^{(x)}(s_t|s_{t-2}) \equiv -\log q_s(s_t, x|s_{t-2}),$$

where $h_z$, $h_x$, and $h_s$ correspond to transition probabilities of $Z$, $X$, and $S$, respectively; $h_x^{(s)}$ and $h_s^{(x)}$ correspond to the transition probabilities of $X$ and $S$ via detailed channels, respectively. These optimal code lengths quantify the transition information with respect to $Z$, $X$, and $S$, respectively, with bits. They can be taken as the dynamical observables associating with the system transitions.

According to the expressions of the transition probabilities $q$, $q_x$, and $q_s$, we note that the observables of the transitions ($h$) can be constructed from the observables of the controls ($l$). This is due to the interactions between the subsystems. More explicitly, we have

$$h_z(x_t, s_t|x_{t-1}, s_{t-1}) = l_x(x_t|s_{t-1}) + l_s(s_t|x_{t-1}),$$
$$h_x(x_t|x_{t-2}) = \log \sum_s 2^{h_x^{(s)}(x_t|x_{t-2})}$$
$$= \log \sum_s 2^{l_x(x_t|s)+l_s(s|x_{t-2})},$$
$$h_s(s_t|s_{t-2}) = \log \sum_x 2^{h_s^{(x)}(s_t|s_{t-2})}$$
$$= \log \sum_x 2^{l_s(s_t|x)+l_x(x|s_{t-2})},$$
$$h_x^{(s)}(x_t|x_{t-2}) = l_x(x_t|s) + l_s(s|x_{t-2}),$$
$$h_s^{(x)}(s_t|s_{t-2}) = l_s(s_t|x) + l_x(x|s_{t-2}).$$

Thus, $h_z$ contains the detailed information of interaction that $S$ controls $X$ [$l_x(x_t|s_{t-1})$] and $X$ controls $S$ [$l_s(s_t|x_{t-1})$] at a transition from $t-1$ to $t$. $h_x^{(s)}$ provides the detailed information that $X$ controls $S$ [$l_s(s|x_{t-2})$] at $t-2$ and $S$ feeds back to $X$ at $t$ [$l_x(x_t|s)$]. $h_s^{(x)}$ measures the detailed information that $S$ controls $X$ [$l_x(x|s_{t-2})$] at $t-2$ and $X$ feeds back to $S$ at $t$ [$l_s(s_t|x)$]. $h_x$ and $h_s$ average the detailed information carried by $h_x^{(s)}$ and $h_s^{(x)}$, respectively, by taking all possible channels

into account. Then, it would be appropriate to measure the information that one subsystem gains at specified state when it controls a state transition of the other subsystem from $t-2$ to $t$ by using the differences

$$h_x^{(s)}(x_t|x_{t-2}) - h_x(x_t|x_{t-2})$$
$$= -\log \frac{q_x(x_t, s|x_{t-2})}{q_x(x_t|x_{t-2})} = -\log q_x(s|x_{t-2}, x_t),$$
$$h_s^{(x)}(s_t|s_{t-2}) - h_s(s_t|s_{t-2})$$
$$= -\log \frac{q_s(s_t, x|s_{t-2})}{q_s(s_t|s_{t-2})} = -\log q_s(x|s_{t-2}, s_t).$$

Here, the $q_x(s|x_{t-2}, x_t) = \frac{q_x(x_t, s|x_{t-2})}{q_x(x_t|x_{t-2})}$ and $q_s(x|s_{t-2}, s_t) = \frac{q_s(s_t, x|s_{t-2})}{q_s(s_t|s_{t-2})}$ represent the probabilities that the specified states $s$ and $x$ conditioning on the transitions $x_{t-2} \to x_t$ and $s_{t-2} \to s_t$, respectively. This can be seen directly from the chain rule of probabilities. Then, $h_x^{(s)}(x_t|x_{t-2}) - h_x(x_t|x_{t-2})$ measures the information gain of $S$ at $s$ when it controls the transition $x_{t-2} \to x_t$; $h_s^{(x)}(s_t|s_{t-2}) - h_s(s_t|s_{t-2})$ measures the information gain of $X$ at $x$ when it controls the transition $s_{t-2} \to s_t$.

### C. Associate dynamical observables with entropy production rates and time-irreversible mutual information rate

By noting the meaning of the dynamical observables of transitions, we use the differences between these observables of backward and forward transitions to define the net bits of information gain or loss of the systems at every transition, respectively, as follows:

$$dh_z(x_t, s_t|x_{t-1}, s_{t-1}) \equiv h_z(x_{t-1}, s_{t-1}|x_t, s_t)$$
$$- h_z(x_t, s_t|x_{t-1}, s_{t-1}) = \log \frac{q(x_t, s_t|x_{t-1}, s_{t-1})}{q(x_{t-1}, s_{t-1}|x_t, s_t)},$$
$$dh_x(x_t|x_{t-2}) \equiv h_x(x_{t-2}|x_t) - h_x(x_t|x_{t-2}) = \log \frac{q_x(x_t|x_{t-2})}{q_x(x_{t-2}|x_t)},$$
$$dh_s(s_t|s_{t-2}) \equiv h_s(s_{t-2}|s_t) - h_s(s_t|s_{t-2}) = \log \frac{q_s(s_t|s_{t-2})}{q_s(s_{t-2}|s_t)}.$$

Here, the corresponding system ($X$, $S$, or $Z$) loses net bits of information ($dh > 0$) or gains net bits of information ($dh < 0$) via the interactions with the environments at transition from $t-m$ to $t$ ($m$ is exact memory length of the system). By taking the averages of $dh$ over all the possible transitions of the corresponding systems, we have

$$\langle dh_z \rangle = R(Z) \geqslant 0, \quad \langle dh_x \rangle = R(X) \geqslant 0, \quad \text{and} \quad \langle dh_s \rangle = R(S) \geqslant 0.$$

This means the information of the system under concern is always dissipated into the environments irreversibly with the rate of bits measured by the EPR.

We should emphasize that the meanings of the EPRs of interacting (non-Markovian) information systems are quite different from those of the systems which are driven by time-invariant environments because one system always behaves as the time-variant environment of the other system. This can be seen from the relations among the EPRs $R(Z)$, $R(X)$, and $R(S)$ directly in this model. The observables with respect to transitions of subsystems via channels $h_x^{(s)}$ and $h_s^{(x)}$ build the bridge among the EPRs. We use the differences between these observables of backward and forward transitions via the same

channel to define the net bits of information gain or loss of the systems at every transition as follows:

$$dh_x^{(s)}(x_t|x_{t-2}) \equiv h_x^{(s)}(x_{t-2}|x_t) - h_x^{(s)}(x_t|x_{t-2})$$

$$= \log \frac{q_x(x_t, s|x_{t-2})}{q_x(x_{t-2}, s|x_t)},$$

$$dh_s^{(x)}(s_t|s_{t-2}) \equiv h_s^{(x)}(s_{t-2}|s_t) - h_s^{(x)}(s_t|s_{t-2})$$

$$= \log \frac{q_s(s_t, x|s_{t-2})}{q_s(s_{t-2}, x|s_t)}.$$

$dh_x^{(s)}$ is recognized as the net bits of information gain ($dh_x^{(s)} < 0$) or loss ($dh_x^{(s)} > 0$) of system $X$ at transitions $x_{t-2} \rightarrow x_t$ via the "channel" $s$; and $dh_s^{(x)}$ is recognized as the net bits of information gain or loss of system $S$ at transitions $s_{t-2} \rightarrow s_t$ via the "channel" $x$. By averaging $dh_x^{(s)}$ and $dh_s^{(x)}$ over all the possible transitions of the corresponding channels, we have the sub-EPRs of subsystems corresponding to channels,

$$R^{(s)}(X) \equiv \langle dh_x^{(s)} \rangle, \quad R^{(x)}(S) \equiv \langle dh_s^{(x)} \rangle.$$

Clearly, the EPR of the composite system $Z$ can be recast by series of these sub-EPRs,

$$R(Z) = \sum_s R^{(s)}(X) = \sum_x R^{(x)}(S).$$

As it is shown in Sec. V B, the differences

$$dh_x^{(s)}(x_t|x_{t-2}) - dh_x(x_t|x_{t-2}) = \log \frac{q_x(s|x_{t-2}, x_t)}{q_x(s|x_t, x_{t-2})},$$

$$dh_s^{(x)}(s_t|s_{t-2}) - dh_s(s_t|s_{t-2}) = \log \frac{q_s(x|s_{t-2}, s_t)}{q_s(x|s_t, s_{t-2})}$$

measure the net bits of information that one system gains or loses when it controls the other system. By taking the averages of these differences over all the possible transitions and channels of the corresponding systems, we have

$$\langle dh_x^{(s)} - dh_x \rangle = R(S|X), \quad \langle dh_s^{(x)} - dh_s \rangle = R(X|S).$$

In this model, the explanation of $R(S|X)$ is that it measures the information dissipation rate of system $S$ when it works as the time-variant environment of $X$. $R(X|S)$ measures the information dissipation rate of system $X$ when it works as the time-variant environment of $S$.

Equation (11) [$R(X|S) = 2I_B(X, S) + R(X)$ and $R(S|X) = 2I_B(X, S) + R(S)$] indicates that one part of $R(S|X)$ [or $R(X|S)$] supplies the self-information dissipation $R(S)$ [or $R(X)$], the other part maintains the time-irreversible MIR $I_B(X, S)$ between systems in this model.

### D. Bounds of time-irreversible part of mutual information rate

By using the FMA method, the EPRs $R(X)$ and $R(S)$ have lower bounds $R(\mathfrak{X}^{(1)})$ and $R(\mathfrak{S}^{(1)})$, respectively, where $\mathfrak{X}^{(1)}$ and $\mathfrak{S}^{(1)}$ refer to the Markov approximations of $X$ and $S$, respectively. The corresponding transition probabilities can be given by

$$f_x(x_t|x_{t-1}) = \sum_{s_{t-1}} \pi(s_{t-1}|x_{t-1})\epsilon(x_t|s_{t-1}) = \pi_x(x_t),$$

$$f_s(s_t|s_{t-1}) = \sum_{x_{t-1}} \pi(x_{t-1}|s_{t-1})d(s_t|x_{t-1}) = \pi_s(s_t),$$

where $\pi(s_{t-1}|x_{t-1}) = \frac{\pi(x_{t-1}, s_{t-1})}{\pi_x(x_{t-1})}$ and $\pi(x_{t-1}|s_{t-1}) = \frac{\pi(x_{t-1}, s_{t-1})}{\pi_s(s_{t-1})}$ are the stationary conditional probabilities with respect to $S$ and $X$, respectively. These two conditional probabilities work in the Markov approximations like $d(s_{t-1}|x_{t-2})$ and $\epsilon(x_{t-1}|s_{t-2})$ in the exact dynamics [$q_x(x_t|x_{t-2})$ and $q_s(s_t|s_{t-2})$], respectively. However, the transition information at time transitions $t - 2 \rightarrow t - 1$ and $t - 1 \rightarrow t$ are totally lost both in $f_x$ and $f_s$. The stationary distributions of the Markov approximations are thus $\pi_x$ and $\pi_s$, respectively. The corresponding approximations of the EPRs $R(X)$ and $R(S)$ can be given by

$$R(\mathfrak{X}^{(1)}) = 0, \quad R(\mathfrak{S}^{(1)}) = 0.$$

Consequentially, we have the approximated time-irreversible MIR

$$I_B(\mathfrak{X}^{(1)}, \mathfrak{S}^{(1)}) = \tfrac{1}{2}R(Z).$$

By noting Eq. (14), we then have the lower and upper bounds of the time-irreversible MIR $I_B(X, S)$ that

$$\tfrac{1}{2}\max\{R(S), R(X)\} \leqslant I_B(X, S) \leqslant \tfrac{1}{2}R(Z).$$

## VII. CONCLUSION

In this work, we quantify the nonequilibrium information dynamics of the two interacting non-Markovian systems: the Markov-type master equation with memory dependence depicts the self-evolution; the nonequilibrium non-Markovian information dynamics of the composite system and mutual information rate depicts the interactions. We characterize the time irreversibility of non-Markovian information dynamics by applying landscape-flux theory. The key point of this theory is the decomposition of the driving force in a Markov-type master equation into time-reversible part (detailed balance preserving landscape part) and time-irreversible part (detailed balance breaking nonequilibrium flux part). Correspondingly, the time-irreversible part of mutual information rate turns out to be closely related to the entropy production rates of subsystems and composite system. Based on our study, we propose the finite memory approximation method which can be used to analyze the time irreversibility or the nonequilibriumness of non-Markovian processes explicitly. The proposed example and the corresponding analysis show the validity of our method.

## ACKNOWLEDGMENT

## APPENDIX A

Let $C$ be a finite-state, discrete-time, irreducible, ergodic, and stationary system with memory length of $m$. Let $\mathcal{C}$ be its state space. Let $\pi$ be its stationary joint distribution of state sequences $v$ with length of $m$. Let $q(v_t|v_{t-1})$ be the transition probability. Then, $C$ is time irreversible beyond $m$ if and only if the stationary flux $J(v_t|v_{t-1})$ is vanishing. That is to say, the

probabilities

$$\Pr(c_1, c_2, \ldots, c_T) = \Pr(c_T, c_{T-1}, \ldots, c_1)$$

hold for $T > m$ iff $J(v_t | v_{t-1}) = 0$.

To prove this, we just need to show that the conclusion holds for $T = m + 2$. We have

$$\Pr(c_1, \ldots, c_{m+2}) - \Pr(c_{m+2}, \ldots, c_1)$$
$$= \pi(v_m) q(v_{m+1}|v_m) q(v_{m+2}|v_{m+1})$$
$$- \pi(\widetilde{v}_{m+2}) q(\widetilde{v}_{m+1}|\widetilde{v}_{m+2}) q(\widetilde{v}_m|\widetilde{v}_{m+1}). \quad (A1)$$

By noting Eq. (3), we proceed with (A1) and have

$$\Pr(c_1, \ldots, c_{m+2}) - \Pr(c_{m+2}, \ldots, c_1)$$
$$= \pi(v_m)(D_1 + B_1)(D_2 + B_2)$$
$$- \pi(\widetilde{v}_{m+2})(\widetilde{D}_2 + \widetilde{B}_2)(\widetilde{D}_1 + \widetilde{B}_1), \quad (A2)$$

where

$$D_1 = D(v_{m+1}|v_m), \quad B_1 = B(v_{m+1}|v_m),$$
$$D_2 = D(v_{m+2}|v_{m+1}), \quad B_2 = B(v_{m+2}|v_{m+1}),$$
$$\widetilde{D}_2 = D(\widetilde{v}_{m+1}|\widetilde{v}_{m+2}), \quad \widetilde{B}_2 = B(\widetilde{v}_{m+1}|\widetilde{v}_{m+2}),$$
$$\widetilde{D}_1 = D(\widetilde{v}_m|\widetilde{v}_{m+1}), \quad \widetilde{B}_1 = B(\widetilde{v}_m|\widetilde{v}_{m+1}).$$

By noting that

$$\pi(v_{t-1}) D(v_t|v_{t-1}) = \pi(\widetilde{v}_t) D(\widetilde{v}_{t-1}|\widetilde{v}_t),$$
$$\pi(v_{t-1}) B(v_t|v_{t-1}) = -\pi(\widetilde{v}_t) B(\widetilde{v}_{t-1}|\widetilde{v}_t),$$

we have $\Pr(c_1, \ldots, c_{m+2}) - \Pr(c_{m+2}, \ldots, c_1) = 0$ iff $B_1 = B_2 = 0$. We can complete the proof for arbitrary $T > m$ by using mathematical induction.

## APPENDIX B

Here is a geometrical interpretation for $J_x$. If $X$ is Markovian ($m = 1$), the stationary joint distribution of its two-state sequences $\Pr(x_{t-1}, x_t)$ forms a matrix $[\Pr(x_{t-1}, x_t)]$. Then, its probability flux $J_x(x_t|x_{t-1})$ is the antisymmetrical part of $\Pr(x_{t-1}, x_t)$ which also forms a matrix (a second-order tensor), namely, $[J_x(x_t|x_{t-1})] = [\Pr(x_{t-1}, x_t)] - [\Pr(x_{t-1}, x_t)]^\dagger = [\Pr(x_{t-1}, x_t)] - [\Pr(x_t, x_{t-1})]$. Here, the dagger symbol represents the matrix transpose. If $X$ is non-Markovian and it has a memory length of $m > 1$, then the joint stationary distribution of $(m + 1)$ sequences $\Pr(\chi_{t-1}, \chi_t) = \pi_x(\chi_{t-1}) q_x(\chi_t|\chi_{t-1})$ [where $(\chi_{t-1}, \chi_t) \equiv [x_{t-m}, \ldots, x_t]$] forms a tensor of order $m + 1$, denoted by $[\Pr(\chi_{t-1}, \chi_t)]$, defined on a coordinates system $[x_{t-m}, \ldots, x_t]$.

To obtain its "tensor"-transpose $[\Pr(\chi_{t-1}, \chi_t)]^\dagger$, we fix all other coordinates except $x_{t-m+k}$ and $x_{t-k}$ [$k = 0, \ldots, (m-1)/2$] and obtain a series of matrices of $(x_{t-m+k}, x_{t-k})$ at all fixed coordinates, namely, the series of matrices $[\Pr(\ldots, x_{t-m+k}, \ldots, x_{t-k}, \ldots)]$. We then transpose these matrices and obtain $[\Pr(\ldots, x_{t-m+k}, \ldots, x_{t-k}, \ldots)]^\dagger = [\Pr(\ldots, x_{t-k}, \ldots, x_{t-m+k}, \ldots)]$. We do these transpose operations from $k = 0$ to $k = (m-1)/2$ to guarantee that all the matrices of the coordinates $(x_{t-m+k}, x_{t-k})$ have been transposed.
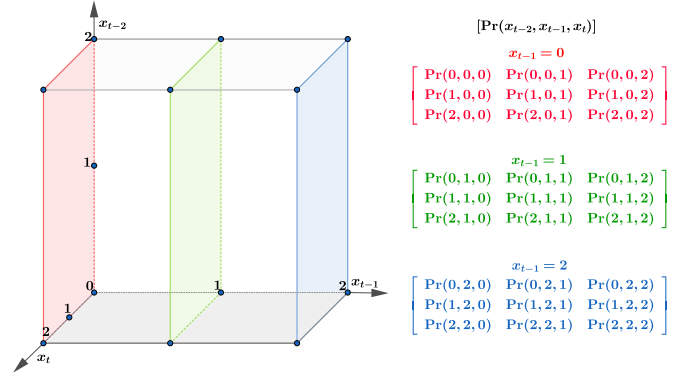


FIG. 3. Geometrical interpretation for the probability flux of a non-Markovian system $X$ with memory. $X$ has memory length of 2 with state space $\mathcal{X} = \{0, 1, 2\}$. The joint probabilities $\Pr(\chi_{t-1}, \chi_t) = \Pr(x_{t-2}, x_{t-1}, x_t)$ are arranged into a tensor space of order 3, i.e., into the lattices of a cube. By fixing $x_{t-1} = 0, 1, 2$, we obtain a series of matrices $[\Pr(x_{t-2}, 0, x_t)]$ (red), $\Pr(x_{t-2}, 1, x_t)]$ (green), and $[\Pr(x_{t-2}, 2, x_t)]$ (blue), with different colors corresponding to the planes shown in the cube. The probability flux which is also a third-order tensor is obtained by calculating the antisymmetrical part of each matrix, i.e., $[J_x] = [\Pr(\chi_{t-1}, \chi_t)] - [\Pr(\chi_{t-1}, \chi_t)]^\dagger = [\Pr(\chi_{t-1}, \chi_t)] - [\Pr(\widetilde{\chi}_t, \widetilde{\chi}_{t-1})]$.

Then, we obtain the "tensor" transpose $[\Pr(\chi_{t-1}, \chi_t)]^\dagger = [\Pr(\widetilde{\chi}_t, \widetilde{\chi}_{t-1})]$ where $\Pr(\widetilde{\chi}_t, \widetilde{\chi}_{t-1}) = \pi_x(\widetilde{\chi}_t) q_x(\widetilde{\chi}_{t-1}|\widetilde{\chi}_t)$. Similar to Markovian case, we have $[J_x(\chi_t|\chi_{t-1})] = [\Pr(\chi_{t-1}, \chi_t)] - [\Pr(\widetilde{\chi}_t, \widetilde{\chi}_{t-1})]$ to be the the antisymmetrical part of tensor $[\Pr(\chi_{t-1}, \chi_t)]$.

A case where $X$ has memory length of 2 has been shown in Fig. 3. It then can be easily verified that $[J_x(\chi_t|\chi_{t-1})]$ is an antisymmetrical tensor of order $m + 1$ such that $[J_x(\chi_t|\chi_{t-1})]^\dagger = [J_x(\widetilde{\chi}_{t-1}|\widetilde{\chi}_t)] = -[J_x(\chi_t|\chi_{t-1})]$.

## APPENDIX C

Here, we derive the exact form of mutual information rate [MIR, Eq. (6)] and the entropy production rate [EPR, Eq. (8)] in steady state. At first, we introduce two related quantities: the "forward" entropy rate $H(C)$ and "backward" entropy rate $\widetilde{H}(C)$ of process $C$, where the description of $C$ has been given in Appendix A. These two quantities are defined by

$$H(C) = \lim_{T \to \infty} \frac{1}{T} \langle -\log \Pr(\Gamma(T)) \rangle_{\Gamma(T)}, \quad (C1)$$

where $\Gamma(T)$ is arbitrary possible time sequence of $C$ in time $T$, and

$$\widetilde{H}(C) = \lim_{T \to \infty} \frac{1}{T} \langle -\log \Pr(\widetilde{\Gamma}(T)) \rangle_{\Gamma(T)}, \quad (C2)$$

where $\widetilde{\Gamma}(T)$ is the time reversal of $\Gamma(T)$.

The MIR [Eq. (14)] can be rewritten into the combination of forward entropy rates as

$$I(X, S) = H(X) + H(S) - H(X, S). \quad (C3)$$

The EPR [Eq. (27)] can be rewritten into the difference between the backward entropy rate and forward entropy rate as

$$R(C) = \widetilde{H}(C) - H(C), \quad \text{for} \quad C = X, S, \text{ or } (X, S). \quad (C4)$$

Thus, we can obtain the MIR and EPR by evaluating $H(C)$ and $\widetilde{H}(C)$. In information theory, these two entropy rates can be obtained from a typical sequence of $C$, where "typical" means in $\Gamma(T)$

(1) the number of the occurrences of a state sequence $\nu$ in $\Gamma(T)$ is

$$\pi(\nu)T + o(T), \text{ for large } T; \tag{C5}$$

(2) the number of the transitions from a memory $\nu_t$ to $\nu_{t+1}$ in $\Gamma(T)$ is

$$\pi(\nu_{t-1})q(\nu_t|\nu_{t-1})T + o(T), \text{ for large } T. \tag{C6}$$

Then, according to the law of large numbers, we have

$$H(C) = -\lim_{T\to\infty} \frac{1}{T} \log \Pr(\Gamma(T)) = -\lim_{T\to\infty} \frac{1}{T} \left\{ \frac{\sum_{\nu_t} \pi(\nu_t) \log \pi(\nu_t)}{+T \sum_{\nu_{t-1}} \sum_{\nu_t} \pi(\nu_{t-1})q(\nu_t|\nu_{t-1}) \log q(\nu_t|\nu_{t-1}) + o(T)} \right\}$$

$$= -\sum_{\nu_t} \sum_{\nu_{t-1}} \pi(\nu_{t-1})q(\nu_t|\nu_{t-1}) \log q(\nu_t|\nu_{t-1}). \tag{C7}$$

Similarly, we have

$$\widetilde{H}(C) = -\sum_{\nu_t} \sum_{\nu_{t-1}} \pi(\nu_{t-1})q(\nu_t|\nu_{t-1}) \log q(\widetilde{\nu}_{t-1}|\widetilde{\nu}_t). \tag{C8}$$

By substituting (C7) and (C8) into (C3) and (C4), respectively, we then have Eqs. (6) and (8).

## APPENDIX D

The log-sum inequality shows that two positive functions $f$ and $g$ with respect to two variables $\eta$ and $\psi$ satisfy the following inequality:

$$\sum_{\eta,\psi} f(\eta, \psi) \log \frac{f(\eta, \psi)}{g(\eta, \psi)} \geqslant \sum_{\eta} \left( \sum_{\psi} f(\eta, \psi) \right) \log \frac{\sum_{\psi} f(\eta, \psi)}{\sum_{\psi} g(\eta, \psi)},$$

where the equality holds if and only if $f(\eta, \psi)/g(\eta, \psi)$ is not a function of $\psi$.

The definition of the EPR suggests that

$$R(X, S) = \lim_{T\to\infty} \sum_{\Gamma_s(T),\Gamma_x(T)} \Pr(\Gamma_x(T), \Gamma_s(T)) \log \frac{\Pr(\Gamma_x(T), \Gamma_s(T))}{\Pr(\widetilde{\Gamma}_x(T), \widetilde{\Gamma}_s(T))}$$

$$\geqslant \lim_{T\to\infty} \sum_{\Gamma_x(T)} \left( \sum_{\Gamma_s(T)} \Pr(\Gamma_x(T), \Gamma_s(T)) \right) \log \frac{\sum_{\Gamma_s(T)} \Pr(\Gamma_x(T), \Gamma_s(T))}{\sum_{\Gamma_s(T)} \Pr(\Gamma_x(T), \Gamma_s(T))} = R(X).$$

Also, the EPRs of two approximations with two memory lengths $m_1$ and $m_2$ ($m_1 > m_2$) satisfy

$$R(\mathfrak{X}, m_1) = \sum_{x_1,\dots,x_{m_1+1}} \Pr(x_1, \dots, x_{m_1+1}) \log \frac{\Pr(x_1, \dots, x_{m_1+1})}{\Pr(x_{m_1+1}, \dots, x_1)}$$

$$\geqslant \sum_{x_1,\dots,x_{m_2+1}} \left( \sum_{x_{m_1+1},\dots,x_{m_2+1}} \Pr(x_1, \dots, x_{m_1+1}) \right) \log \frac{\sum_{x_{m_1+1},\dots,x_{m_2+1}} \Pr(x_1, \dots, x_{m_1+1})}{\sum_{x_{m_1+1},\dots,x_{m_2+1}} \Pr(x_{m_1+1}, \dots, x_1)} = R(\mathfrak{X}, m_2).$$

where we use the substitutions of variables $\eta = [x_1, ..., x_{m_2+1}]$ and $\psi = [x_{m_1+1}, ..., x_{m_2+1}]$ for applying the log-sum inequality.

## APPENDIX E

Assume that $\pi$ is the unique stationary distribution of the composite system $Z$. This means that the transition matrix $[q(z_t|z_{t-1})]$ has a unique eigenvalue of 1 and $\pi$ is the corresponding eigenvector which is a probability distribution. Thus, $\pi$ satisfies that

$$\pi(x_t, s_t) = \sum_{x_{t-1},s_{t-1}} \pi(x_{t-1}, s_{t-1})q(x_t, s_t|x_{t-1}, s_{t-1}). $$

Let $\pi_x(x) = \sum_s \pi(x, s)$ and $\pi_s(s) = \sum_x \pi(x, s)$ be the two marginal stationary distributions. We can easily verify that

$$\sum_{s_{t-1}} \pi_s(s_{t-1})\epsilon(x_t|s_{t-1}) = \sum_{x_t} \sum_{x_{t-1},s_{t-1}} \pi(x_{t-1}, s_{t-1})q(x_t, s_t|x_{t-1}, s_{t-1}) = \pi_x(x_t),$$

$$\sum_{x_{t-1}} \pi_x(x_{t-1})d(s_t|x_{t-1}) = \sum_{s_t} \sum_{x_{t-1},s_{t-1}} \pi(x_{t-1}, s_{t-1})q(x_t, s_t|x_{t-1}, s_{t-1}) = \pi_s(s_t).$$

Let $\rho(x, s) = \pi_x(x)\pi_s(s)$ be the direct product distribution. We then have

$$\sum_{x_{t-1},s_{t-1}} \rho(x_{t-1}, s_{t-1})q(x_t, s_t|x_{t-1}, s_{t-1}) = \sum_{x_{t-1},s_{t-1}} \rho(x_{t-1}, s_{t-1})\epsilon(s_{t-1}|x_t)d(s_t|x_{t-1})$$

$$= \left(\sum_{s_{t-1}} \pi_s(s_{t-1})\epsilon(x_{t-1}|s_t)\right)\left(\sum_{x_{t-1}} \pi_x(x_{t-1})d(s_t|x_{t-1})\right) = \rho(x_t, s_t).$$

Thus, $\rho$ is a stationary distribution of $Z$. Since $\pi$ is the unique stationary distribution of $Z$, then we must have $\pi = \rho$.

The FMA method shows that, for $m \geqslant 2$,

$$q_x\left(\chi_t^{(m)}|\chi_{t-1}^{(m)}\right) = \frac{\sum_{s_1,...,s_t} \Pr(x_1, s_1)\epsilon(x_2|s_1)d(s_2|x_1)\ldots\epsilon(x_t|s_{t-1})d(s_t|x_{t-1})}{\sum_{s_1,...,s_{t-1}} \Pr(x_1, s_1)\epsilon(x_2|s_1)d(s_2|x_1)\ldots\epsilon(x_{t-1}|s_{t-2})d(s_{t-1}|x_{t-2})}$$

$$= \frac{\left(\sum_{s_1} \Pr(x_1, s_1)\epsilon(x_2|s_1)\right)\left(\sum_{s_2} d(s_2|x_1)\epsilon(x_3|s_2)\right)\ldots\left(\sum_{s_{t-1}} d(s_{t-1}|x_{t-2})\epsilon(x_t|s_{t-1})\right)}{\left(\sum_{s_1} \Pr(x_1, s_1)\epsilon(x_2|s_1)\right)\left(\sum_{s_2} d(s_2|x_1)\epsilon(x_3|s_2)\right)\ldots\left(\sum_{s_{t-2}} d(s_{t-2}|x_{t-3})\epsilon(x_{t-1}|s_{t-2})\right)}$$

$$= \sum_{s_{t-1}} d(s_{t-1}|x_{t-2})\epsilon(x_t|s_{t-1}) = q_x\left(\chi_t^{(2)}|\chi_{t-1}^{(2)}\right) = q_x(x_t|x_{t-2}).$$

Similarly, we have $q_s(\varsigma_t^{(m)}|\varsigma_{t-1}^{(m)}) = \sum_{x_{t-1}} \epsilon(x_{t-1}|s_{t-2}) d(s_t|x_{t-1})$ for $m \geqslant 2$.

To justify the numerical results of FMA, we can use a conventional method which evaluates $R(X)$, $R(S)$, $I(X, S)$, and $I_B(X, S)$ directly from a typical sequence of $Z$ (see [15]). For large time $T$, the corresponding results can be given by

$$R(X) \approx \frac{1}{T} \log \frac{\Pr(\Gamma_x(T))}{\Pr(\widetilde{\Gamma}_x(T))}, \quad R(S) \approx \frac{1}{T} \log \frac{\Pr(\Gamma_s(T))}{\Pr(\widetilde{\Gamma}_s(T))},$$

$$I(X, S) \approx \frac{1}{T} \log \frac{\Pr(\Gamma_z(T))}{\Pr(\Pr(\Gamma_x(T))\Pr(\Gamma_s(T)))},$$

$$I_B(X, S) \approx \frac{1}{2T}\left(\log \frac{\Pr(\Gamma_z(T))}{\Pr(\widetilde{\Gamma}_z(T))} - \log \frac{\Pr(\Gamma_x(T))}{\Pr(\widetilde{\Gamma}_x(T))} - \log \frac{\Pr(\Gamma_s(T))}{\Pr(\widetilde{\Gamma}_s(T))}\right),$$

where $\Gamma_z(T) = (\Gamma_x(T), \Gamma_s(T))$ is a typical sequence of $Z$ [hence, $\Gamma_x(T)$ and $\Gamma_s(T)$ are typical sequences of $X$ and $S$, respectively] in time $T$; $\widetilde{\Gamma}_z(T)$, $\widetilde{\Gamma}_x(T)$, and $\widetilde{\Gamma}_s(T)$ are the corresponding time-reversal sequences. The convergence of this method can be observed as $T$ increases.

For numerical simulations, we simply let the state space of the system $X$ be $\mathcal{X} = \{x : x = 1, 2, 3, 4\}$ and the state space of the system $S$ be $\mathcal{S} = \{s : s = 1, 2, 3\}$. The conditional probabilities $\epsilon$ read as

$$\epsilon = \begin{bmatrix} 0.1168 & 0.5300 & 0.0591 \\ 0.2437 & 0.0093 & 0.3049 \\ 0.3219 & 0.3808 & 0.4136 \\ 0.3175 & 0.0799 & 0.2224 \end{bmatrix}.$$

The conditional probabilities $d$ read as

$$d = \begin{bmatrix} 0.4525 & 0.4445 & 0.4231 & 0.3962 \\ 0.4539 & 0.2110 & 0.1870 & 0.1812 \\ 0.0935 & 0.3446 & 0.3899 & 0.4226 \end{bmatrix}.$$

Here, values of $\epsilon$ and $d$ are arranged into the matrices of $[\epsilon(x_t|s_{t-1})]$ (with rows being labeled by $x$ and columns being labeled by $s$) and $[d(s_t|x_{t-1})]$ (with rows being labeled by $s$ and columns being labeled by $x$), respectively. Consequentially, the exact transition probabilities of $X$ and $S$ read as

$$q_x = \begin{bmatrix} 0.2990 & 0.1841 & 0.1716 & 0.1673 \\ 0.1430 & 0.2153 & 0.2237 & 0.2271 \\ 0.3572 & 0.3659 & 0.3687 & 0.3713 \\ 0.2008 & 0.2346 & 0.2360 & 0.2343 \end{bmatrix},$$

$$q_s = \begin{bmatrix} 0.4232 & 0.4367 & 0.4254 \\ 0.2222 & 0.3283 & 0.2088 \\ 0.3546 & 0.2350 & 0.3658 \end{bmatrix}.$$

Here, values of $q_x$ and $q_s$ are arranged into the matrices of $[q_x(x_t|x_{t-2})]$ and $[q_s(s_t|s_{t-2})]$ with rows being labeled by states at $t$ and columns being labeled by states at $t - 2$, respectively. The exact stationary probabilities of $X$ and $S$

TABLE II. Numerical results of $R(Z)$, $R(X)$, $R(S)$, $I(X, S)$, $I_B(X, S)$, $I_B^L(X, S)$, and $I_B^U(X, S)$.

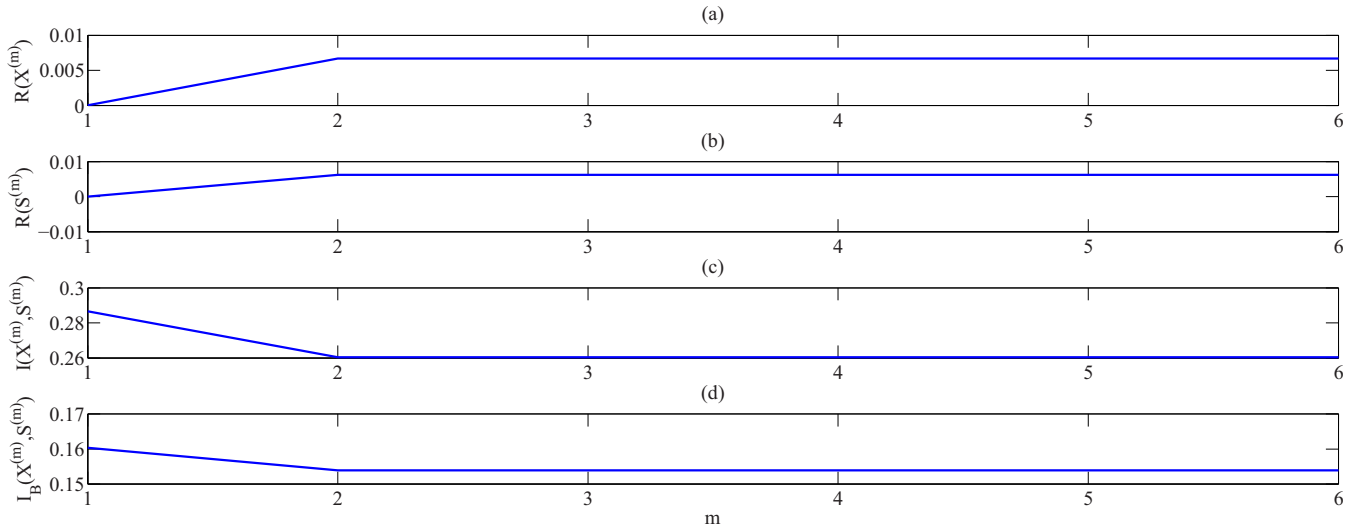| $R(Z)$ | $R(X)$ | $R(S)$ | $I(X, S)$ | $I_B(X, S)$ | $I_B^L(X, S)$ | $I_B^U(X, S)$ |
|---|---|---|---|---|---|---|
| 0.3208 | 0.0067 | 0.0062 | 0.2603 | 0.1540 | −0.0031 | 0.1604 |

FIG. 4. The values of (a) $R(\mathfrak{X}^{(m)})$, (b) $R(\mathfrak{S}^{(m)})$, (c) $I(\mathfrak{X}^{(m)}, \mathfrak{S}^{(m)})$, and (d) $I_B(\mathfrak{X}^{(m)}, \mathfrak{S}^{(m)})$ from $m = 1$ to 6.

read as

$$\pi_x = [\pi_x(1), \pi_x(2), \pi_x(3), \pi_x(4)]$$
$$= [0.1985, 0.2067, 0.3665, 0.2283],$$
$$\pi_s = [\pi_s(1), \pi_s(2), \pi_s(3)] = [0.4272, 0.2436, 0.3292].$$

We evaluate $R(\mathfrak{X}^{(m)})$, $R(\mathfrak{S}^{(m)})$, $I(\mathfrak{X}^{(m)}, \mathfrak{S}^{(m)})$, and $I_B(\mathfrak{X}^{(m)}, \mathfrak{S}^{(m)})$ by using FMA. We select the thresholds of relative error $\delta = 10^{-5}$ for both entities. All the calculations are terminated at $M = 2$. This demonstrates the conclusion that the processes of the subsystems are both

non-Markovian chains with memory lengths of 2. In fact, we have evaluated the model with unified memory lengths from $m = 1$ to 6 to check whether there exists any exception (see Fig. 4). We also calculate the lower and upper bounds of $I_B(X, S)$, $I_B^L(X, S) = \max\{-\frac{1}{2}R(X), -\frac{1}{2}R(S)\}$, and $I_B^U(X, S) = \frac{1}{2}R(Z)$, respectively. The values of numerical results are listed in Table II.

We also plot the curves of $R(X)$, $R(S)$, $I(X, S)$, and $I_B(X, S)$ with increasing $T$ ($1 \leqslant T \leqslant 10^6$) by using the conventional method as the comparison to the FMA method (see Fig. 5).
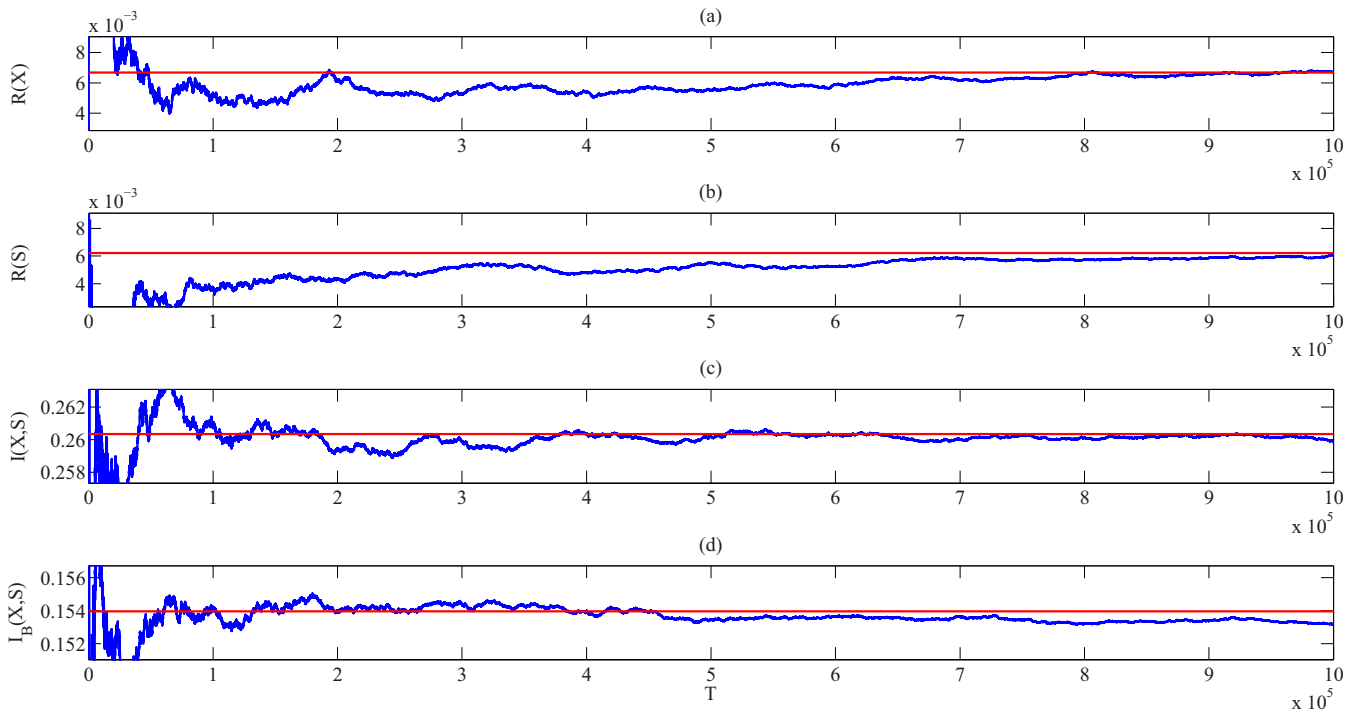


FIG. 5. Comparisons of (a) $R(X)$, (b) $R(S)$, (c) $I(X, S)$, and (d) $I_B(X, S)$ by using conventional method and exact values. Curved lines, conventional method. Horizontal lines, exact values with $m = 2$.

[1] P. Strasberg, G. Schaller, T. Brandes, and M. Esposito, Phys. Rev. Lett. **110**, 040601 (2013).

[2] J. V. Koski, A. Kutvonen, I. M. Khaymovich, T. Ala-Nissila, and J. P. Pekola, Phys. Rev. Lett. **115**, 260602 (2015).

[3] T. Mcgrath, N. S. Jones, P. R. Ten Wolde, and T. E. Ouldridge, Phys. Rev. Lett. **118**, 028101 (2017).

[4] T. Sagawa and M. Ueda, Phys. Rev. Lett. **109**, 180602 (2012).

[5] J. M. Horowitz and M. Esposito, Phys. Rev. X **4**, 031015 (2014).

[6] J. M. R. Parrondo, J. M. Horowitz, and T. Sagawa, Nat. Phys. **11**, 131 (2015).

[7] B. L. Mark and Y. Ephraim, IEEE Trans. Signal Process. **62**, 2709 (2014).

[8] Y. Ephraim and B. L. Mark, Found. Trends Signal Process. **6**, 1 (2013).

[9] A. C. Barato, D. Hartich, and U. Seifert, J. Stat. Phys. **153**, 460 (2013).

[10] Q. Zeng and J. Wang, Entropy **19**, 678 (2017).

[11] J. Wang, L. Xu, and E. K. Wang, Proc. Natl. Acad. Sci. USA **105**, 12271 (2008).

[12] J. Wang, Adv. Phys. **64**, 1 (2015).

[13] C. H. Li, E. K. Wang, and J. Wang, J. Chem. Phys J. Chem. Phys. **136**, 194108 (2012).

[14] C. E. Shannon, Bell Syst. Tech. J. **27**, 379 (1948).

[15] T. M. Cover and J. A. Thomas, *Elements of Information Theory* (Wiley, Hoboken, NJ, 2003).

[16] U. Seifert, Rep. Prog. Phys. **75**, 126001 (2012).

[17] J. Jacod and P. Protter, *Probability Essentials* (Springer, Berlin, 2000).

[18] S.-J. Wu and M. T. Chu, Appl. Math. Comput. **303**, 226 (2017).

[19] R. Gray and J. Kieffer, IEEE Trans. Inf. Theory **26**, 412 (1980).

[20] C. Maes, F. Redig, and A. van Moffaert, J. Math. Phys. **41**, 1528 (2000).

[21] P. Gaspard, J. Stat. Phys. **117**, 599 (2004).

[22] V. Y. Chernyak, M. Chertkov, and C. Jarzynski, J. Stat. Mech.: Theor. Exp. (2006) P08001.

[23] D. Mandal and C. Jarzynski, Proc. Natl. Acad. Sci. USA **109**, 11641 (2012).

[24] A. C. Barato and U. Seifert, Phys. Rev. Lett. **112**, 090601 (2014).

[25] M. Esposito and G. Schaller, Europhys. Lett. **99**, 30003 (2012).

[26] D. Huffman, Proc. IRE **40**, 1098 (1952).