# Information Landscape and Flux, Mutual Information Rate Decomposition and Connections to Entropy Production

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

[1]  State Key Laboratory of Electroanalytical Chemistry, Changchun Institute of Applied Chemistry, Changchun, Jilin 130022, China; qzeng@ciac.ac.cn
[2]  Department of Chemistry and Physics, State University of New York, Stony Brook, NY 11794, USA
*   Correspondence: jin.wang.1@stonybrook.edu; Tel.: +1-631-632-1185

**Abstract:** We explored the dynamics of two interacting information systems. We show that for the Markovian marginal systems, the driving force for information dynamics is determined by both the information landscape and information flux. While the information landscape can be used to construct the driving force to describe the equilibrium time-reversible information system dynamics, the information flux can be used to describe the nonequilibrium time-irreversible behaviors of the information system dynamics. The information flux explicitly breaks the detailed balance and is a direct measure of the degree of the nonequilibrium or time-irreversibility. We further demonstrate that the mutual information rate between the two subsystems can be decomposed into the equilibrium time-reversible and nonequilibrium time-irreversible parts, respectively. This decomposition of the Mutual Information Rate (MIR) corresponds to the information landscape-flux decomposition explicitly when the two subsystems behave as Markov chains. Finally, we uncover the intimate relationship between the nonequilibrium thermodynamics in terms of the entropy production rates and the time-irreversible part of the mutual information rate. We found that this relationship and MIR decomposition still hold for the more general stationary and ergodic cases. We demonstrate the above features with two examples of the bivariate Markov chains.

## 1. Introduction

There is growing interest in studying two interacting information systems in the fields of control theory, information theory, communication theory, nonequilibrium physics and biophysics [1–9]. Significant progresses has been made recently towards the understanding of the information system in terms of information thermodynamics [10–13]. However, the identification of the global driving forces for the information system dynamics is still challenging. Here, we aim to fill this gap by quantifying the driving forces for the information system dynamics. Inspired by the recent development of landscape and flux theory for the continuous nonequilibrium systems [14–16] and the Markov chain decomposition dynamics for the discrete systems [17–23], we show that at least for the underlying marginal Markovian cases, the driving force for information dynamics is determined by both the information landscape and information flux. The information landscape can be used to construct the driving force responsible for the equilibrium time-reversible part of the information dynamics. The information flux explicitly breaks the detailed balance and provides a quantitative measure of the degree of nonequilibrium or time-irreversibility. It is responsible for the time-irreversible part of the information dynamics. The Mutual Information Rate (MIR) [24] represents the correlation

between two information subsystems. We uncovered that the MIR between the two subsystems can be decomposed into the time-reversible and time-irreversible parts, respectively. Especially when the two subsystems act as Markov chains, this decomposition can be expressed in terms of information landscape-flux decomposition for Markovian dynamics. An important signature of nonequilibrium is the Entropy Production Rate (EPR) [17,25,26]. We also uncover the intimate relation between the EPRs and the time-irreversible part of the MIR. We demonstrate the above features with two cases of the bivariate Markov chains. Furthermore, we show that the decomposition of the MIR and the relationship between the EPRs and the time-irreversible part of the MIR still hold for more general stationary and ergodic cases.

## 2. Bivariate Markov Chains

Markov chains have been often assumed for the underlying dynamics of the total system in random environments. When the two subsystems together jointly form a Markov chain in continuous or discrete time, the resulting chain is called the *Bivariate Markov Chain* (BMC, a special case of the multivariate Markov chain with two stochastic variables). The processes of the two subsystems are correspondingly said to be marginal processes or a marginal chain. The BMC was used to model ion channel currents [2]. It was also used to model delays and congestion in a computer network [3]. Recently, different models of BMC appeared in nonequilibrium statistical physics for capturing or implementing Maxwell's demon [4–6], which can be seen as one marginal chain in the BMC playing feedback control to the other marginal chain. Although the BMC has been studied for decades, there are still challenges on quantifying the dynamics of the whole, as well as the two subsystems. This is because neither of them needs to be a Markovian chain in general [7], and the quantifications of the probabilities (densities) for the trajectories of the two subsystems involve the complicated random matrix multiplications [8]. This leads to the problem not exactly being analytically solvable. The corresponding numerical solutions often lack direct mathematical and physical interpretations.

The conventional analysis of the BMC focuses on the mutual information [9] of the two subsystems for quantifying the underlying information correlations. There are three main representations of this. The first one was proposed and emphasized in the works of Sagawa, T. and Ueda, M. [11] and Parrondo, J. M. R., Horowitz, J. M. and Sagawa, T. [10], respectively, for explaining the mechanism of Maxwell's demon in Szilard's engine. In this representation, the mutual information between the demon and controlled system characterizes the observation and the feedback of the demon. This leads to an elegant approach, which includes the increment of the mutual information into a unified fluctuation relation. The second representation was proposed by the work of Horowitz, J. M. and Esposito, M. [12] in an attempt to explain the violation of the second law in a specified BMC, the bipartite model, where the mutual information is divided into two parts corresponding to the two subsystems, respectively, which were said to be the information flows. This representation tries to explain the mechanism of the demon because one can see that the information flows do contribute to the entropy production for both the demon and controlled system. The first two representations are based on the ensembles of the subsystem states. This means that the mutual information is defined only on the time-sliced distributions of the system states, which somehow lack the information of subsystem dynamics: the time-correlations of the observation and feedback of the demon. The last representation was seen in the work of Barato, A. C., Hartich, D. and Seifert, U. [13], where a more general definition of mutual information in information theory was used, which is defined on the trajectories of the two subsystems. More exactly, this is the so-called *Mutual Information Rate* (MIR) [24], which quantifies the correlation between the two subsystem dynamics. However, due to the difficulties from the possible underlying non-Markovian property of the marginal chains, exactly solvable models and comprehensive conclusions are still challenging from this representation.

In this study, we study the discrete-time BMC in both stochastic information dynamics and thermodynamics. To avoid the technical difficulty caused by non-Markovian dynamics, we first assume that the two marginal chains follow the Markovian dynamics. The non-Markovian case will

be discussed elsewhere. We explore the time-irreversibility of BMC and marginal processes in the steady state. Then, we decompose the driving force for the underlying dynamics as the information landscape and information flux [14–16], which can be used to describe the time-reversible parts and time-irreversible parts, respectively. We also prove that the non-vanishing flux fully describes the time-irreversibility of BMC and marginal processes.

We focus on the mutual information rate between the two marginal chains. Since the two marginal chains are assumed to be Markov chains here, the mutual information rate is exactly analytically solvable, which can be seen as the averaged conditional correlation between the two subsystem states. Here, the conditional correlations reveal the time correlations between the past states and the future states.

Corresponding to the landscape-flux decomposition in stochastic dynamics, we decompose the MIR into two parts: the time-reversible and time-irreversible parts, respectively. The time-reversible part measures the part of the correlations between the two marginal chains in both forward and backward processes of BMC. The time-irreversible part measures the difference between the correlations in forward and backward processes of BMC, respectively. We can see that a non-vanishing time-irreversible part of the MIR must be driven by a non-vanishing flux in the steady state, and it can be seen as the sufficient condition for a BMC to be time-irreversible.

We also reveal the important fact that the time-irreversible parts of MIR contribute to the nonequilibrium *Entropy Production Rate* (EPR) of the BMC by the simple equality:

*EPR of BMC = EPR of 1st marginal chain + EPR of 2nd marginal chain + 2 × time-irreversible part of MIR.*

The decomposition of the MIR and the relation between the time-irreversible part of MIR and EPRs can also be found in stationary and ergodic non-Markovian cases, which will be given in the discussions in the Appendix. This may help to develop a general theory of nonequilibrium non-Markovian interacting information systems.

## 3. Information Landscape and Information Flux for Determining the Information Dynamics, Time-Irreversibility

Consider the case that two interacting information systems form a finite-state, discrete-time, ergodic and irreducible bivariate Markov chain,

$$Z = (X, S) = \{(X(t), S(t)), t \geq 1\}, \tag{1}$$

We assume that the information state space of $X$ is given by $\mathcal{X} = \{1, ..., d\}$ and the information state space of $S$ is given by $\mathcal{S} = \{1, ..., l\}$. The information state space of $Z$ is then given by $\mathcal{Z} = \mathcal{X} \times \mathcal{S}$. The stochastic information dynamics can then be quantitatively described by the time evolution of the probability distribution of information state space $Z$, characterized by the following master equation (or the information system dynamics) in discrete time,

$$p_z(z; t+1) = \sum_{z'} q_z(z|z') p_z(z'; t), \text{ for } t \geq 1, \text{ and } z \in \mathcal{Z} \tag{2}$$

where $p_z(z; t) = p_z(x, s; t)$ is the probability of observing state $z$ (or joint probability of $X = x$ and $S = s$) at time $t$; $q_z(z|z') = q_z(x, s|x', s') \geq 0$ are the transition probabilities from $z' = (x', s')$ to $z = (x, s)$, respectively, and have $\sum_z q_z(z|z') = 1$.

We assume that there exists a unique stationary distribution $\pi_z$ such that $\pi_z(z) = \sum_{z'} q_z(z|z') \pi_z(z')$. Then, given an arbitrary initial probability distribution, the probability distribution goes to $\pi_z$ exponentially fast in time. If the initial distribution is $\pi_z$, we say that $Z$ is in *Steady State* (SS), and our discussion is based on this SS.

The marginal chains of $Z$, i.e., $X$ and $S$, do not need to be Markov chains in general. For the simplicity of analysis, we assume that both marginal chains are Markov chains, and the corresponding

transition probabilities are given by $q_x(x|x')$ and $q_s(s|s')$ (for $x, x' \in \mathcal{X}$ and $s, s' \in \mathcal{S}$), respectively. Then, we have the following master equations (or the information system dynamics) for $X$ and $S$, respectively,

$$p_x(x; t+1) = \sum_{x'} q_x(x|x') p_x(x'; t), \tag{3}$$

and,

$$p_s(s; t+1) = \sum_{s'} q_s(s|s') p_s(s'; t), \tag{4}$$

where $p_x(x; t)$ and $p_s(s; t)$ are the probabilities of observing $X = x$ and $S = s$ at time $t$, respectively.

We consider that both Equations (3) and (4) have unique stationary solutions $\pi_x$ and $\pi_s$, which satisfy $\pi_x(x) = \sum_{x'} q_x(x|x')\pi_x(x')$ and $\pi_s(s) = \sum_{s'} q_s(s|s')\pi_s(s')$ respectively. Furthermore, we assume that when $Z$ is in SS, $\pi_x$ and $\pi_s$ are also achieved. The relations between $\pi_x$, $\pi_s$ and $\pi_z$ read,

$$\begin{cases} \pi_x(x) = \sum_s \pi_z(x, s), \\ \pi_s(s) = \sum_x \pi_z(x, s). \end{cases} \tag{5}$$

In the rest of this paper, we let $X^T = \{X(1), X(2), ..., X(T)\}$, $S^T = \{S(1), S(2), ..., S(T)\}$, and $Z^T = \{Z(1), Z(2), ..., Z(T)\} = (X^T, S^T)$ denote the time sequences of $X$, $S$ and $Z$ in time $T$, respectively.

To characterize the time-irreversibility of the Markov chain $C$ in information dynamics in SS, we introduce the concept of probability flux. Here, we let $C$ denote the arbitrary Markov chain in $\{Z, X, S\}$, and let $c$, $\pi_c$, $q_c$ and $C^T$ denote arbitrary state of $C$, the stationary distribution of $C$, the transition probabilities of $C$ and a time sequence of $C$ in time $T$ and in SS, respectively.

The averaged number transitions from the state $c'$ to state $c$, denoted by $N(c' \to c)$, in unit time in SS can be obtained as:

$$N(c' \to c) = \pi_c(c') q_c(c|c').$$

This is also the probability of the time sequence $C^T = \{C(1) = c', C(2) = c\}$, $(T = 2)$. Correspondingly, the averaged number of reverse transitions, denoted by $N(c \to c')$, reads:

$$N(c \to c') = \pi_c(c) q_c(c'|c).$$

This is also the probability of the time-reverse sequence $\widetilde{C}^T = \{C(1) = c, C(2) = c'\}$, $(T = 2)$. The difference between these two transition numbers measures the time-reversibility of the forward sequence $C^T$ in SS,

$$\begin{aligned} J_c(c' \to c) &= N(c' \to c) - N(c \to c') \\ &= P(C^T) - P(\widetilde{C}^T) \\ &= \pi_c(c') q_c(c|c') - \pi_c(c) q_c(c'|c), \text{ for } C = X, S, \text{ or } Z. \end{aligned} \tag{6}$$

Then, $J_c(c' \to c)$ is said to be the probability flux from $c'$ to $c$ in SS. If $J_c(c' \to c) = 0$ for arbitrary $c'$ and $c$, then $C^T$ $(T = 2)$ is time-reversible; otherwise, when $J_c(c' \to c) \neq 0$, $C^T$ is time-irreversible. Clearly, we have from Equation (6) that:

$$J_c(c' \to c) = -J_c(c \to c'). \tag{7}$$

The transition probability determines the evolution dynamics of the information system. We can decompose the transition probabilities $q_c(c|c')$ into two parts: the time-reversible part $D_c$ and time-irreversible part $B_c$, which read:

$$q_c(c|c') = D_c(c' \to c) + B_c(c' \to c), \text{ with}$$
$$\begin{cases} D_c(c' \to c) = \frac{1}{2\pi_c(c')}(\pi_c(c')q_c(c|c') + \pi_c(c)q_c(c'|c)), \\ B_c(c' \to c) = \frac{1}{2\pi_c(c')}J_c(c' \to c). \end{cases} \tag{8}$$

From this decomposition, we can see that the information system dynamics is determined by two driving forces. One of the driving forces is determined by the steady state probability distribution. This part of the driving force is time-reversible. The other driving force for the information dynamics is the steady state probability flux, which breaks the detailed balance and quantifies the time-irreversibility. Since the steady state probability distribution measures the weight of the information state, therefore it can be used to quantify the *information landscape*. If we define the potential landscape for the information system as $\phi = -\log\pi$, then the driving force $D_c(c' \to c) = \frac{1}{2}(q_c(c|c') + \frac{\pi_c(c)}{\pi_c(c')}q_c(c'|c)) = \frac{1}{2}(q_c(c|c') + \exp[-(\phi_c(c) - \phi_c(c')]q_c(c'|c))$ is expressed in term of the difference of the potential landscape. This is analogous to the landscape-flux decomposition of Langevin dynamics in [15]. Notice that the information landscape is directly related to the steady state probability distribution of the information system. In general, the information landscape is at nonequilibrium since the detailed balance is often broken for general cases. Only when the detailed balance is preserved, the nonequilibrium information landscape is reduced to the equilibrium information landscape. Even though the information landscape is not at equilibrium in general, the driving force $D_c(c' \to c)$ is time-reversible due to the decomposition construction. The steady state probability flux measures the information flow in the dynamics and therefore can be termed as the *information flux*. In fact, the nonzero information flux explicitly breaks the detailed balance because of the net flow to or from the system. It is therefore a direct measure of the degree of the nonequilibrium or time-irreversibility in terms of the detailed balance breaking.

Note that the decomposition for the discrete Markovian information process can be viewed as the separation of the current corresponding to the $2B_c(c' \to c)\pi_c(c')$ here and the activity corresponding to the $2D_c(c' \to c)\pi_c(c')$ in a previous study [19]. The landscape and flux decomposition here for the reduced information dynamics are in a similar spirit as the whole state space decomposition with the information system and the associated environments. When the detailed balance is broken, the information landscape (defined as the negative logarithm of the steady state probability $\phi = -\log\pi$) is not the same as the equilibrium landscape under the detailed balance. There can be uniqueness issue related to the decomposition. To avoid the confusion, we make a physical choice, or in other words, we can fix the gauge so that the information landscape always coincides with the equilibrium landscape when the detailed balance is satisfied. In other words, we want to make sure the Boltzmann law applies at equilibrium with detailed balance. In this way, we can decompose the information landscape and information flux for nonequilibrium information systems without detailed balance. By solving the linear master equation for the steady state, we can quantify the nonequilibrium information landscape, and from that, we can obtain the corresponding steady state probability flux. Some studies discussed various aspects of this issue [18,19,27,28].

By Equations (7) and (8), we have the following relations:

$$\begin{cases} \pi_c(c')D_c(c' \to c) = \pi_c(c)D_c(c \to c'), \\ \pi_c(c')B_c(c' \to c) = -\pi_c(c)B_c(c \to c'). \end{cases} \tag{9}$$

As we can see in the next section, $D_c$ and $B_c$ are useful for us to quantify time-reversible and time-irreversible observables of $C$, respectively.

We give the interpretation that the non-vanishing information flux $J_c$ fully measures the time-irreversibility of the chain $C$ in time $T$ for $T \geq 2$. Let $C^T$ be an arbitrary sequence of $C$ in SS, and without loss of generality, we let $T = 3$. Similar to Equation (6), the measure of the time-irreversibility of $C^T$ can be given by the difference between the probability of $C^T = \{C(1), C(2), C(3)\}$ and that of its time-reversal $\widetilde{C}^T = \{C(3), C(2), C(1)\}$, such as:

$$
\begin{aligned}
&P(C^T) - P(\widetilde{C}^T) \\
&= \pi_c(C(1))q_c(C(2)|C(1))q_c(C(3)|C(2)) - \pi_c(C(3))q_c(C(2)|C(3))q_c(C(1)|C(2)) \\
&= \pi_c(C(1)) \left( D_c(C(1) \to C(2)) + B_c(C(1) \to C(2)) \right) \left( D_c(C(2) \to C(3)) + B_c(C(2) \to C(3)) \right) - \\
&\quad \pi_c(C(3)) \left( D_c(C(3) \to C(2)) + B_c(C(3) \to C(2)) \right) \left( D_c(C(2) \to C(1)) + B_c(C(2) \to C(1)) \right), \\
&\quad \text{for } C = X, S \text{ or } Z.
\end{aligned}
$$

Then, by the relations given in Equation (9), we have that $P(C^T) - P(\widetilde{C}^T) = 0$ holds for arbitrary $C^T$ if and only if $B_c(C(1) \to C(2)) = B_c(C(2) \to C(3)) = 0$ or equivalently $J_c(C(1) \to C(2)) = J_c(C(2) \to C(3)) = 0$. This conclusion can be made for arbitrary $T > 3$. Thus, non-vanishing $J_c$ can fully describe the time-irreversibility of $C$ for $C = X, S$ or $Z$.

We show the relations between the fluxes of the whole system $J_z$ and of the subsystem $J_x$ as follows:

$$
\begin{aligned}
J_x(x' \to x) &= \pi_x(x')q_x(x|x') - \pi_x(x)q_x(x'|x) \\
&= P(\{x', x\}) - P(\{x, x'\}) \\
&= \sum_{s,s'} \left( P(\{(x', s'), (x, s)\}) - P(\{(x, s), (x', s')\}) \right) \\
&= \sum_{s,s'} \left( \pi_z(x', s')q_z(x, s|x', s') - \pi_z(x, s)q_z(x', s'|x, s) \right) \\
&= \sum_{s,s'} J_z((x', s') \to (x, s)).
\end{aligned}
$$

(10)

Similarly, we have:

$$
J_s(s' \to s) = \sum_{x,x'} J_z((x', s') \to (x, s)).
$$

(11)

These relations indicate that the subsystem fluxes $J_x$ and $J_s$ can be seen as the coarse-grained levels of total system flux $J_z$ by averaging over the other parts of the system $S$ and $X$, respectively. We should emphasize that non-vanishing $J_z$ does not mean $X$ or $S$ is time-irreversible and vice versa.

## 4. Mutual Information Decomposition to Time-Reversible and Time-Irreversible Parts

According to information theory, the two interacting information systems represented by bivariate Markov chain $Z$ can be characterized by the *Mutual Information Rate* (MIR) between the marginal chains $X$ and $S$ in SS. The mutual information rates represent the correlation between two interacting information systems. The MIR is defined on the probabilities of all possible time sequences, $P(Z^T)$, $P(X^T)$ and $P(S^T)$ and is given by [24],

$$
I(X, S) = \lim_{T \to \infty} \frac{1}{T} \sum_{Z^T} P(Z^T) \log \frac{P(Z^T)}{P(X^T)P(S^T)}.
$$

(12)

It measures the correlation between $X$ and $S$ in unit time, or say, the efficient bits of information that $X$ and $S$ exchange with each other in unit time. The MIR must be non-negative, and a vanishing $I(X, S)$ indicates that $X$ and $S$ are independent of each other. More explicitly, the corresponding probabilities of these sequences can be evaluated by using Equations (2)–(4); we have:

$$\begin{cases} P(X^T) = \pi_x(X(1)) \prod_{t=1}^{T-1} q_x(X(t+1)|X(t)), \\ P(S^T) = \pi_s(S(1)) \prod_{t=1}^{T-1} q_s(S(t+1)|S(t)), \\ P(Z^T) = \pi_z(Z(1)) \prod_{t=1}^{T-1} q_z(Z(t+1)|Z(t)). \end{cases}$$

By substituting these probabilities into Equation (12) (see Appendix A), we have the exact expression of MIR as:

$$\begin{aligned} I(X,S) &= \sum_{z,z'} \pi_z(z') q_z(z|z') \log \frac{q_z(z|z')}{q_x(x|x') q_s(s|s')} \\ &= \langle i(z|z') \rangle_{z',z} \geq 0, \text{ for } z = (x,s), \text{ and } z' = (x',s'). \end{aligned} \tag{13}$$

where $i(z|z') = \log \frac{q_z(z|z')}{q_x(x|x') q_s(s|s')}$ is the conditional (Markovian) correlation between the states $x$ and $s$ when the transition $z' = (x',s') \to z = (x,s)$ occurs. This indicates that when the two marginal processes are both Markovian, the MIR is the average of the conditional (Markovian) correlations. These correlations are measurable when transitions occur, and they can be seen as the observables of $Z$.

By noting the decomposition of transition probabilities in Equation (8), we have a corresponding decomposition of $I(X,S)$ such as:

$$I(X,S) = I_D(X,S) + I_B(X,S), \text{ with}$$
$$\begin{cases} I_D(X,S) = \sum_{z,z'} \pi_z(z') D_z(z|z') i(z|z') = \frac{1}{2} \sum_{z,z'} (\pi_z(z') q_z(z|z') + \pi_z(z) q_z(z'|z)) i(z|z'), \\ I_B(X,S) = \sum_{z,z'} \pi_z(z') B_z(z|z') i(z|z') = \frac{1}{2} \sum_{z,z'} J_z(z|z') i(z|z') = \frac{1}{4} \sum_{z,z'} J_z(z|z') (i(z|z') - i(z'|z)). \end{cases} \tag{14}$$

This means that the mutual information representing the correlations between the two interacting systems can be decomposed into the time-reversible equilibrium part and the time-irreversible nonequilibrium part. The origin of this is from the fact that the underlying information system dynamics is determined by both the time-reversible information landscape and time-irreversible information flux. These equations are very important to establish the link to the time-irreversibility. We now give further interpretation for $I_D(X,S)$ and $I_B(X,S)$:

Consider a bivariate Markov chain $Z$ in SS wherein $X$ and $S$ are dependent on each other, i.e., $I(X,S) = I_D(X,S) + I_B(X,S) > 0$. By the ergodicity of $Z$, we have the MIR, which measures the averaged conditional correlation along the time sequences $Z^T$,

$$\lim_{T \to \infty} \frac{1}{T} \langle i(Z(t+1)|Z(t)) \rangle_{Z^T} = I(X,S), \text{ for } 1 < t < T.$$

Then, $I_B(X,S)$ measures the change of the averaged conditional correlation between $X$ and $S$ when a sequence of $Z$ turns back in time,

$$\lim_{T \to \infty} \frac{1}{T} \langle i(Z(t+1)|Z(t)) - i(Z(t)|Z(t+1)) \rangle_{Z^T} = 2 I_B(X,S).$$

A negative $I_B(X,S)$ shows that the correlation between $X$ and $S$ becomes strong in the time-reversal process of $Z$; A positive $I_B(X,S)$ shows that the correlation becomes weak in the time-reversal process of $Z$. Both cases show that the $Z$ is time-irreversible since we have a non-vanishing $J_z$. However, the case of $I_B(X,S) = 0$ is complicated, since it indicates either a vanishing $J_z$ or a non-vanishing $J_z$. Anyway, we see that a non-vanishing $I_B(X,S)$ is a sufficient condition for $Z$ to be time-irreversible. On the other hand, $I_D(X,S) = I(X,S) - I_B(X,S)$ measures the correlation remaining in the backward process of $Z$.

The definition of MIR in Equation (12) turns out to be appropriate for even more general stationary and ergodic (Markovian or non-Markovian) processes. Consequentially, the decomposition of MIR is useful to quantify the correlation between two stationary and ergodic processes in a wider sense, i.e., to monitor the changes of the correlation in the forward and the backward processes. As a special case,

the analytical expressions in Equation (14) are the reduced results, which are valid for Markovian cases. A brief discussion of the decomposition of MIR of more general processes can be found in Appendix B.

## 5. Relationship between Mutual Information and Entropy Production

The *Entropy Production Rates* (EPR) or energy dissipation (cost) rate at steady state is a quantitative nonequilibrium measure, which characterizes the time-irreversibility of the underlying processes. The EPR of a stationary and ergodic process $C$ (here $C = Z$, $X$ or $S$) can be given by the difference between the averaged surprisal (negative logarithmic probability) of the backward sequences $\widetilde{C}^T$ and that of forward sequences $C^T$ in the long time limit, i.e.,

$$
\begin{aligned}
R_c &= \lim_{T \to \infty} \frac{1}{T} \left\langle \log P(C^T) - \log P(\widetilde{C}^T) \right\rangle_{C^T} \\
&= \lim_{T \to \infty} \frac{1}{T} \left\langle \log \frac{P(C^T)}{P(\widetilde{C}^T)} \right\rangle_{C^T} \geq 0,
\end{aligned}
\tag{15}
$$

where $R_c$ is said to be the EPR of $C$ [25]; $-\log P(C^T)$ and $-\log P(\widetilde{C}^T)$ are said to be the surprisal of a forward and a backward sequence of $C$, respectively. We see that $C$ is time-reversible (i.e., $P(C^T) = P(\widetilde{C}^T)$ for arbitrary $C^T$ for large $T$) if and only if $R_c = 0$. Additionally, this is due to the form of $R_c$, which is exactly a Kullback–Leibler divergence. When $C$ is Markovian, then $R_c$ reduces into the following form when $Z$, $X$ or $S$ is assigned to $C$, respectively [17,26],

$$
\begin{cases}
R_z = \frac{1}{2} \sum_{z,z'} J_z(z' \to z) \log \frac{q_z(z|z')}{q_z(z'|z)}, \\
R_x = \frac{1}{2} \sum_{x,x'} J_x(x' \to x) \log \frac{q_x(x|x')}{q_x(x'|x)}, \\
R_s = \frac{1}{2} \sum_{s,s'} J_s(s' \to s) \log \frac{q_s(s|s')}{q_s(s'|s)},
\end{cases}
\tag{16}
$$

where total and subsystem entropy productions $R_z$, $R_x$ and $R_s$ correspond to $Z$, $X$ and $S$, respectively. Here, $R_z$ usually contains the detailed interaction information of the system (or subsystems) and environments; $R_x$ and $R_s$ provide the coarse-grained information of time-irreversible observables of $X$ and $Z$, respectively. Each non-vanishing EPR indicates that the corresponding Markov chain is time-irreversible. Again, we emphasize that a non-vanishing $R_z$ does not mean $X$ or $S$ is time-irreversible and vice versa.

We are interested in the connection between these EPRs and mutual information. We can associate them with $I_B(X, S)$ by noting Equations (10), (11) and (14). We have:

$$
\begin{aligned}
I_B(X, S) &= \frac{1}{4} \sum_{z,z'} J_z(z|z')(i(z|z') - i(z'|z)) \\
&= \frac{1}{4} \sum_{z,z'} J_z(z|z') \log \frac{q_z(z|z')}{q_z(z'|z)} - \frac{1}{4} \sum_{x,x'} J_x(x|x') \log \frac{q_x(x|x')}{q_x(x'|x)} - \frac{1}{4} \sum_{s,s'} J_s(s|s') \log \frac{q_s(s|s')}{q_s(s'|s)} \\
&= \frac{1}{2}(R_z - R_x - R_s).
\end{aligned}
\tag{17}
$$

We note that $I_B(X, S)$ is intimately related to the EPRs. This builds up a bridge between these EPRs and the irreversible part of the mutual information. Moreover, we also have:

$$
\begin{cases}
R_z = R_x + R_s + 2I_B(X, S) \geq 0, \\
R_x + R_s \geq -2I_B(X, S), \\
R_z \geq 2I_B(X, S).
\end{cases}
\tag{18}
$$

This indicates that the time-irreversible MIR contributes to the detailed EPRs. In other words, the differences of the entropy production rate of the whole system and subsystems provide the origin of the time-irreversible part of the mutual information. This reveals the nonequilibrium thermodynamic

origin of the irreversible mutual information or correlations. Of course, since the EPR is related to the flux directly as is seen from the above definitions, the origin of the EPR or nonequilibrium thermodynamics is from the non-vanishing information flux for the nonequilibrium dynamics. On the other hand, the irreversible part of the mutual information measures the correlations, and it contributes to the EPRs of the correlated subsystems.

Furthermore, the last expression in Equation (17) (also the expressions in Equation (18)) can be generalized to more general stationary and ergodic processes. A related discussion and demonstration of this can be seen in Appendix B.

## 6. A Simple Case: The Blind Demon

As a concrete example, we consider a two-state system coupled to two information baths $a$ and $b$. The states of the system are denoted by $\mathcal{X} = \{x : x = 0, 1\}$, respectively. Each bath sends an instruction to the system. If the system adopts one of them, it then follows the instruction and makes the change of the state. The instructions generated from one bath are independently and identically distributed (Bernoulli trials). Both the probability distributions of the instructions corresponding to the baths follow Bernoulli distributions and read $\{\epsilon_a(x) : x \in \mathcal{X}, \epsilon_a(x) \geq 0, \sum_x \epsilon_a(x) = 1\}$ for bath $a$ and $\{\epsilon_b(x) : x \in \mathcal{X}, \epsilon_b(x) \geq 0, \sum_x \epsilon_b(x) = 1\}$ for bath $b$, respectively. Since the system cannot execute two instructions simultaneously, there exists an information demon that makes choices for the system. The demon is blind to caring about the system, and it makes choices independently and identically distributed. The choices of the demon are denoted by $\mathcal{S} = \{s : s = a, b\}$, respectively. The probability distribution of the demon's choices reads $\{P(s) : s \in \mathcal{S}, P(a) = p, P(b) = 1 - p, p \in [0, 1]\}$. Still, we use $Z = (X, S)$ with $X \in \mathcal{X}$ and $S \in \mathcal{S}$ to denote the BMC of the system and the demon.

Consequentially, the transition probabilities of the system read:

$$q_x(x|x') = p\epsilon_a(x) + (1 - p)\epsilon_b(x).$$

The transition probabilities of the demon read:

$$q_s(s|s') = P(s).$$

Additionally, the transition probabilities of the joint chain read:

$$q_z(x, s|x', s') = P(s)\epsilon_{s'}(x).$$

We have the corresponding steady state distributions or the information landscapes as,

$$\begin{cases} \pi_x(x) = p\epsilon_a(x) + (1 - p)\epsilon_b(x), \\ \pi_s(s) = P(s), \\ \pi_z(x, s) = P(s)\pi_x(x). \end{cases}$$

We obtain the information fluxes as,

$$\begin{cases} J_x(x' \to x) = 0, \text{ for all } x, x' \in \mathcal{X} \\ J_s(s' \to s) = 0, \text{ for all } s, s' \in \mathcal{S} \\ J_z((x', s') \to (x, s)) = P(s)P(s')(\pi_x(x')\epsilon_{s'}(x) - \pi_x(x)\epsilon_s(x')). \end{cases}$$

Here, we use the notations $\epsilon_s(x')$ and $\epsilon_{s'}(x)$ ($s, s' = a$ or $b$) to denote the probabilities of the instructions $x'$ or $x$ from bath $a$ or $b$ briefly. We obtain the EPRs as:

$$\begin{cases} R_x = 0, \\ R_s = 0, \\ R_z = \sum_x p(1-p)(\epsilon_a(x) - \epsilon_b(x))(\log \epsilon_a(x) - \log \epsilon_b(x)). \end{cases}$$

We evaluate the MIR as:

$$I(X,S) = -\sum_x \pi_x(x) \log \pi_x(x) + p \sum_x \epsilon_a(x) \log \epsilon_a(x) + (1-p) \sum_x \epsilon_b(x) \log \epsilon_b(x).$$

The time-irreversible part of $I(X,S)$ reads,

$$I_B(X,S) = \frac{1}{2} R_z.$$

## 7. Conclusions

In this work, we identify the driving forces for the information system dynamics. We show that for marginal Markovian information systems, the information dynamics is determined by both the information landscape and information flux. While the information landscape can be used to construct the driving force for describing the time-reversible behavior of the information dynamics, the information flux can be used to describe the time-irreversible behavior of the information dynamics. The information flux explicitly breaks the detailed balance and provides a quantitative measure of the degree of the nonequilibrium or time-irreversibility. We further demonstrate that the mutual information rate, which represents the correlations, can be decomposed into the time-reversible part and the time-irreversible part originated from the landscape and flux decomposition of the information dynamics. Finally, we uncover the intimate relationship between the difference of the entropy productions of the whole system and those of the subsystems and the time-irreversible part of the mutual information. This will help with understanding the non-equilibrium behavior of the interacting information system dynamics in stochastic environments. Furthermore, we verify that our conclusions on the mutual information rate and entropy production rate decomposition can be made more general for the stationary and ergodic processes.

**Author Contributions:** Qian Zeng and Jin Wang conceived and designed the experiments; Qian Zeng performed the experiments; Qian Zeng and Jin Wang analyzed the data; Qian Zeng and Jin Wang contributed reagents/materials/analysis tools; Qian Zeng and Jin Wang wrote the paper.

**Conflicts of Interest:** The authors declare no conflict of interest.

## Abbreviations

The following abbreviations are used in this manuscript:

BMC    Bivariate Markov Chain
EPR    Entropy Production Rate
MIR    Mutual Information Rate
SS     Steady State

## Appendix A

Here, we derive the exact form of the Mutual Information Rate (MIR, Equation (13)) in the steady state by using the cumulant-generating function.

We write an arbitrary time sequence of $Z$ in time $T$ in the following form:

$$Z^T = \{Z(1), ..., Z(i), ..., Z(T)\}, \text{ for } T \geq 2,$$

where $Z(i)$ (for $i \geq 1$) denotes the state at time $i$. The corresponding probability of $Z^T$ is in the following form:

$$P(Z^T) = \pi_z(Z_1) \left\{ \prod_{i=1}^{T-1} q_z(Z_{i+1}|Z_i) \right\}. \tag{A1}$$

We let the chain $U = (X, S)$ denote a process that $X$ and $S$ follow the same Markov dynamics in $Z$, but are independent of each other. Then, we have that the transition probabilities of $U$ read:

$$q_u(u|u') = q(x,s|x',s') = q_x(x|x')q_s(s|s'). \tag{A2}$$

Then, the probability of a time sequence of $U$, $U^T$, with the same trajectory of $Z^T$ reads:

$$P(U^T) = \pi_u(Z_1) \left\{ \prod_{i=1}^{T-1} q_u(Z_{i+1}|Z_i) \right\}, \tag{A3}$$

with $\pi_u(x,s) = \pi_x(x)\pi_s(s)$ being the stationary probability of $U$.

For evaluating the exact form of MIR, we introduce the cumulant-generating function of the random variable $\log \frac{P(Z^T)}{P(U^T)}$,

$$K(m,T) = \log \left\langle \exp \left( m \log \frac{P(Z^T)}{P(U^T)} \right) \right\rangle_{Z^T}. \tag{A4}$$

We can see that:

$$\begin{aligned}
\lim_{T\to\infty} \lim_{m\to 0} \frac{1}{T} \frac{\partial K(m,T)}{\partial m} \\
= \lim_{T\to\infty} \frac{1}{T} \left\langle \log \frac{P(Z^T)}{P(U^T)} \right\rangle_{Z^T} \\
= I(X,S).
\end{aligned} \tag{A5}$$

Thus, our idea is to evaluate $K(m,T)$ at first. We have:

$$\begin{aligned}
K(m,T) &= \log \left\langle \exp \left( m \log \frac{P(Z^T)}{P(U^T)} \right) \right\rangle_{Z^T} \\
&= \log \left\{ \sum_{Z^T} \frac{(P(Z^T))^{m+1}}{(P(U^T))^m} \right\} \\
&= \log \left\{ \sum_{\{Z(0),Z(1),...,Z(T)\}} \frac{(\pi_z^{m+1}(Z_0))}{(\pi_u^m(Z_0))} \prod_{i=0}^{T-1} \frac{q_z^{m+1}(Z_{i+1}|Z_i)}{q_u^m(Z_{i+1}|Z_i)} \right\},
\end{aligned} \tag{A6}$$

where we realize that the last equality can be rewritten in the form of matrix multiplication.

We introduce the following matrices and vectors for Equation (A6) such that:

$$\begin{aligned}
\mathbf{Q}_z &= \left\{ (\mathbf{Q_z})_{(z,z')} = q_z(z|z'), \text{ for } z,z' \in \mathcal{Z} \right\}, \\
\mathbf{G}(m) &= \left\{ (\mathbf{G}(m))_{(z,z')} = \frac{q_z^{m+1}(z|z')}{q_u^m(z|z')}, \text{ for } z,z' \in \mathcal{Z} \right\}, \\
\boldsymbol{\pi}_z &= \left\{ (\boldsymbol{\pi}_z)_z = \pi_z(z), \text{ for } z \in \mathcal{Z} \right\}, \\
\boldsymbol{v}(m) &= \left\{ (\boldsymbol{v}(m))_z = \frac{\pi_z^{m+1}(z)}{\pi_u^m(z)} \right\},
\end{aligned} \tag{A7}$$

where $\mathbf{Q}_z$ is the transition matrix of $Z$; $\boldsymbol{\pi}_z$ is the stationary distribution of $Z$. It can be also verified that:

$$\boldsymbol{Q}_z = \boldsymbol{G}(0),$$
$$\boldsymbol{\pi}_z = \boldsymbol{v}(0),$$
$$\boldsymbol{\pi}_z = \boldsymbol{Q}_z\boldsymbol{\pi}_z,$$
$$\mathbf{1}^\dagger\boldsymbol{Q}_z = \mathbf{1}^\dagger,$$
$$\lim_{m\to 0}\frac{d\boldsymbol{G}(m)}{dm} = \left\{\left(\lim_{m\to 0}\frac{d\boldsymbol{G}(m)}{dm}\right)_{(z,z')} = q_z(z|z')\log\frac{q_z(z|z')}{q_u(z|z')}, \text{ for } z,z' \in \mathcal{Z}\right\},$$
$$\lim_{m\to 0}\frac{d\boldsymbol{v}(m)}{dm} = \left\{\left(\lim_{m\to 0}\frac{d\boldsymbol{v}(m)}{dm}\right)_z = \pi_z(z)\log\frac{\pi_z(z)}{\pi_u(z)}, \text{ for } z \in \mathcal{Z}\right\},$$

$$(A8)$$

where $\mathbf{1}^\dagger$ is the vector of all ones with appropriate dimension.

Then, $K(m,T)$ can be rewritten in a compact form such that:

$$K(m,T) = \log\left\{\mathbf{1}^\dagger\boldsymbol{G}^{T-1}(m)\boldsymbol{v}(m)\right\}. \tag{A9}$$

Then, we substitute Equation (A9) into Equation (A5) and have:

$$
\begin{aligned}
I(X,S) &= \lim_{T\to\infty}\lim_{m\to 0}\frac{1}{T}\frac{\partial K(m,T)}{\partial m} \\
&= \lim_{T\to\infty}\lim_{m\to 0}\frac{1}{T}\frac{\partial\log\left\{\mathbf{1}^\dagger\boldsymbol{G}^{T-1}(m)\boldsymbol{v}(m)\right\}}{\partial m} \\
&= \lim_{T\to\infty}\lim_{m\to 0}\frac{1}{T}\left\{(T-1)\mathbf{1}^\dagger\boldsymbol{G}^{T-2}(m)\frac{d\boldsymbol{G}(m)}{dm}\boldsymbol{v}(m) + \mathbf{1}^\dagger\boldsymbol{G}^{T-1}(m)\frac{d\boldsymbol{v}(m)}{dm}\right\} \\
&= \lim_{T\to\infty}\frac{1}{T}\left\{(T-1)\mathbf{1}^\dagger\boldsymbol{G}^{T-2}(0)\left(\lim_{m\to 0}\frac{d\boldsymbol{G}(m)}{dm}\right)\boldsymbol{v}(0) + \mathbf{1}^\dagger\boldsymbol{G}^{T-1}(0)\left(\lim_{m\to 0}\frac{d\boldsymbol{v}(m)}{dm}\right)\right\}.
\end{aligned}
$$

$$(A10)$$

By noting Equation (A8) and $T \geq 2$, we obtain Equation (13) from Equation (A10) such that:

$$
\begin{aligned}
I(X,S) &= \lim_{T\to\infty}\frac{1}{T}\left\{(T-1)\mathbf{1}^\dagger\boldsymbol{G}^{T-2}(0)\left(\lim_{m\to 0}\frac{d\boldsymbol{G}(m)}{dm}\right)\boldsymbol{v}(0) + \mathbf{1}^\dagger\boldsymbol{G}^{T-1}(0)\left(\lim_{m\to 0}\frac{d\boldsymbol{v}(m)}{dm}\right)\right\} \\
&= \lim_{T\to\infty}\left\{\left(1-\frac{1}{T}\right)\mathbf{1}^\dagger\left(\lim_{m\to 0}\frac{d\boldsymbol{G}(m)}{dm}\right)\boldsymbol{\pi}_z + \frac{1}{T}\mathbf{1}^\dagger\left(\lim_{m\to 0}\frac{d\boldsymbol{v}(m)}{dm}\right)\right\} \\
&= \mathbf{1}^\dagger\left(\lim_{m\to 0}\frac{d\boldsymbol{G}(m)}{dm}\right)\boldsymbol{\pi}_z \\
&= \sum_{(x,s),(x',s')}\pi_z(x',s')q_z(x,s|x',s')\log\frac{q_z(x,s|x',s')}{q_x(x|x')q_s(s|s')}.
\end{aligned}
$$

$$(A11)$$

## Appendix B

*Appendix B.1 Discussions on the Generality of Mutual Information Rate Decomposition and Connections to Entropy Production in Terms of Equations (14), (17), and (18)*

For general cases, indeed, we do not expect that both $X$ and $S$ are Markovian. Even the joint chain $Z$ may be non-Markovian. This means that Equation (2) may fail to depict the dynamics of $Z$. Then, the landscape-flux decomposition needs to be generalized to this situation. Such decomposition was not developed yet for the non-Markovian cases. This will be discussed in a separate work. However, when $Z$ is a stationary and ergodic process (also assume that both $X$ and $S$ are stationary and ergodic), we show that the MIR can be decomposed into two parts as is shown in Equation (14), and an interesting relation between the MIR and EPRs can still be found in the same form of the last expression in Equation (17).

We are interested in the correlation between the forward sequences of $X$ and $S$, which can be measured by $\log \frac{P(Z^T)}{P(X^T)P(S^T)}$ ($Z^T = (X^T, S^T)$), then the MIR can be used to quantify the average rate of this correlation in the long time limit as shown in Equation (12). Furthermore, we are interested in the averaged difference between the rate of the correlation of the backward processes and that of the forward processes. This comes to the time-irreversible part of the MIR defined by:

$$I_B(X, S) = \lim_{T \to \infty} \frac{1}{2T} \left\langle \log \frac{P(Z^T)}{P(X^T)P(S^T)} - \log \frac{P(\widetilde{Z}^T)}{P(\widetilde{X}^T)P(\widetilde{S}^T)} \right\rangle_{Z^T}, \tag{A12}$$

where $\log \frac{P(\widetilde{Z}^T)}{P(\widetilde{X}^T)P(\widetilde{S}^T)}$ quantifies the correlation between the backward sequences of $X$ and $S$. Clearly, the time-irreversible part of MIR depicting the correlation of the forward processes of $X$ and $S$ is enhanced ($I_B(X, S) > 0$) or weakened ($I_B(X, S) < 0$) compared to that of the backward processes. The other important part of the MIR, namely the time-reversible part, shows that the averaged rate of the correlation that remains in both forward and backward processes,

$$I_D(X, S) = \lim_{T \to \infty} \frac{1}{2T} \left\langle \log \frac{P(Z^T)}{P(X^T)P(S^T)} + \log \frac{P(\widetilde{Z}^T)}{P(\widetilde{X}^T)P(\widetilde{S}^T)} \right\rangle_{Z^T}, \tag{A13}$$

Consequentially, the MIR $I(X, S)$ is decomposed into two parts shown as $I(X, S) = I_D(X, S) + I_B(X, S)$. In Markovian cases, each part of the MIR reduces to the form in Equation (14) respectively.

The relation between the time-irreversible part of the MIR and EPRs can be shown as follows,

$$\begin{aligned} I_B(X, S) &= \lim_{T \to \infty} \frac{1}{2T} \left\langle \log \frac{P(Z^T)}{P(X^T)P(S^T)} - \log \frac{P(\widetilde{Z}^T)}{P(\widetilde{X}^T)P(\widetilde{S}^T)} \right\rangle_{Z^T} \\ &= \lim_{T \to \infty} \frac{1}{2T} \left\{ \left\langle \log \frac{P(Z^T)}{P(\widetilde{Z}^T)} \right\rangle_{Z^T} - \left\langle \log \frac{P(X^T)}{P(\widetilde{X}^T)} \right\rangle_{X^T} - \left\langle \log \frac{P(S^T)}{P(\widetilde{S}^T)} \right\rangle_{S^T} \right\} \\ &= \frac{1}{2} (R_z - R_x - R_s), \end{aligned} \tag{A14}$$

which is in the same form as Equation (17). Additionally, due to the non-negativity of the EPRs, the inequalities in (18) still hold for general cases.

*Appendix B.2. The Smart Demon*

To verify the conclusions in more general cases, we constructed a model of a smart demon, which reflects a more general situation in the nature: the two information subsystems play feedback to each other. A three-state information system is connected to two information baths labeled $a$ and $b$, respectively. The states of the system are denoted by $\mathcal{X} = \{x : x = 0, 1, 2\}$, respectively. Each bath sends an instruction to the system. If the system adopts one of them, it then follows the instruction and makes a change of the state. The instructions generated from one arbitrary bath are independent and identically distributed. The probability distributions of the instructions corresponding to the baths read $\{\epsilon_s(x) : \epsilon_s(x) \geq 0, \sum_{x \in \mathcal{X}} \epsilon_s(x) = 1\}$ (for $s = a, b$), respectively. Since the system cannot execute the two incoming instructions simultaneously, there exists an information demon making choices for the system. The choices of the demon are denoted by the labels of the baths $\mathcal{S} = \{s : s = a, b\}$, respectively. The demon observes the state of the system and plays feedback. The (conditional) probability distribution of the demon's choices reads $\{d(s|x', s') : d(s|x', s') \geq 0, \sum_{s \in \mathcal{S}} d(s|x', s') = 1, x' \in \mathcal{X}, s' \in \mathcal{S}\}$. Still, we use $X$, $S$ and $Z = (X, S)$ to denote the processes of the system, the demon and the corresponding joint chain, a BMC, respectively.

The transition probabilities of the BMC read:

$$q_z(z|z') = q_z(x, s|x', s') = d(s|x', s')\epsilon_s(x),$$

where $\epsilon_s(x)$ denotes the probability of the instruction $x$ from bath $s = a, b$. We assume that there is a unique stationary distribution of $z$, $\pi_z$ such that:

$$\pi_z(z) = \sum_{z'} q_z(z|z') \pi_z(z').$$

The stationary distributions of $S$ and $X$ then read:

$$\begin{cases} \pi_s(s) = \sum_x \pi_z(x, s), \\ \pi_x(x) = \sum_s \pi_z(x, s). \end{cases}$$

The behavior of the demon can be seen as a Markovian process in the steady state. The corresponding transition probabilities of the system read:

$$q_s(s|s') = \frac{1}{\pi_s(s')} \sum_{x'} d(s|x', s') \pi_z(x', s').$$

It can be verified that $\pi_s$ is the unique stationary distribution of $S$. However, the dynamics of the system always behaves as a non-Markovian process in general.

To characterize the time-irreversibility of $Z$, $X$ and $S$, we use the definition of EPR in Equation (15) and have:

$$\begin{cases} R_z = \frac{1}{2} \sum_{z,z'} J_z(z' \to z) \log \frac{q_z(z|z')}{q_z(z'|z)}, \\ R_s = \frac{1}{2} \sum_{s,s'} J_s(s' \to s) \log \frac{q_s(s|s')}{q_s(s'|s)} = 0, \\ R_x = \lim_{T \to \infty} \frac{1}{T} \sum_{X^T} P(X^T) \log \frac{P(X^T)}{P(\tilde{X}^T)}, \end{cases}$$

where:

$$P(X^T) = \sum_{S^T} P(Z^T = (X^T, S^T)).$$

To quantify the correlation between the system and demon, we use the definition of MIR in Equation (12).

We are also interested in the time-irreversible part of MIR, $I_B(X, S)$, which influences the EPR of the system, $R_x$. This can be seen from Equation (A14) such that:

$$R_x = R_z - R_s - 2I_B(X, S).$$

We use numerical simulations, which evaluate $R_x$, $I(X, S)$ and $I_B(X, S)$ directly from the typical sequences of $Z$ (see [7,8]). The corresponding results can be given by:

$$\begin{cases} R_x \approx \frac{1}{T} \log \frac{P(X^T)}{P(\tilde{X}^T)}, \text{ for large } T, \\ I(X, S) \approx \frac{1}{T} \log \frac{P(Z^T)}{P(X^T)P(S^T)}, \text{ for large } T, \\ I_B(X, S) \approx \frac{1}{2T} \log \frac{P(Z^T)}{P(X^T)P(S^T)} - \frac{1}{2T} \log \frac{P(\tilde{Z}^T)}{P(\tilde{X}^T)P(\tilde{S}^T)}, \text{ for large } T, \end{cases}$$

where $Z^T = (X^T, S^T)$ is a typical sequence of $Z$ (hence, $X^T$ and $S^T$ are typical sequences of $X$ and $S$, respectively). The convergence of this numerical simulation can be observed as $T$ increases. To confirm the result $R_x = R_z - R_s - 2I_B(X, S)$, we use different typical sequences in calculating $R_x$ and $I_B(X, S)$, respectively. $R_z$ and $R_s$ are calculated by using the corresponding analytical results shown above.

For numerical simulations, we randomly choose two groups of the parameters: the probabilities of the instructions of the baths $\epsilon_a$ and $\epsilon_b$ and the probabilities of the demon's choices $d$ (see Tables A1 and A2).

We evaluate $R_x$, $I(X, S)$ and $I_B(X, S)$ for both groups. The values of the numerical results are listed in Table A3.

**Table A1.** Two groups of $\epsilon_a$ and $\epsilon_b$.

|  | $\{\epsilon_a(x=0), \epsilon_a(x=1), \epsilon_a(x=2)\}$ | $\{\epsilon_b(x=0), \epsilon_b(x=1), \epsilon_b(x=2)\}$ |
|---|---|---|
| Group 1 | $\{0.2344, 0.2730, 0.4926\}$ | $\{0.4217, 0.4094, 0.1689\}$ |
| Group 2 | $\{0.1305, 0.3972, 0.4723\}$ | $\{0.3358, 0.0010, 0.6633\}$ |

**Table A2.** Two groups of $d$.

|  | $\{d(s=a|x=0, s=a), d(s=b|x=0, s=a)\}$ | $\{d(s=a|x=1, s=b), d(s=b|x=0, s=b)\}$ |
|---|---|---|
| Group 1 | $\{0.3844, 0.6156\}$ | $\{0.6811, 0.3189\}$ |
| Group 2 | $\{0.1072, 0.8928\}$ | $\{0.7473, 0.2527\}$ |

|  | $\{d(s=a|x=1, s=a), d(s=b|x=1, s=a)\}$ | $\{d(s=a|x=1, s=b), d(s=b|x=1, s=b)\}$ |
|---|---|---|
| Group 1 | $\{0.5195, 0.4805\}$ | $\{0.8088, 0.1912\}$ |
| Group 2 | $\{0.6595, 0.3405\}$ | $\{0.1600, 0.8400\}$ |

|  | $\{d(s=a|x=2, s=a), d(s=b|x=2, s=a)\}$ | $\{d(s=a|x=2, s=b), d(s=b|x=2, s=b)\}$ |
|---|---|---|
| Group 1 | $\{0.3775, 0.6225\}$ | $\{0.3340, 0.6660\}$ |
| Group 2 | $\{0.0232, 0.9768\}$ | $\{0.0814, 0.9186\}$ |

**Table A3.** Numerical results of $R_z$, $R_x$, $I(X, S)$ and $I_B(X, S)$.

|  | $R_z$ | $R_x$ | $I(X,S)$ | $I_B(X,S)$ |
|---|---|---|---|---|
| Group 1 | 0.0645 | 0.0018 | 0.0885 | 0.0313 |
| Group 2 | 0.5485 | 0.1291 | 0.3385 | 0.2097 |

## References

1. Shannon, C.E. A mathematical theory of communication. *Bell Syst. Tech. J.* **1948**, *27*, 379–423.
2. Ball, F.; Yeo, G.F. Lumpability and Marginalisability for Continuous-Time Markov Chains. *J. Appl. Probab.* **1993**, *30*, 518–528.
3. Wei, W.; Wang, B.; Towsley, D. Continuous-time hidden Markov models for network performance evaluation. *Perform. Eval.* **2002**, *49*, 129–146.
4. Strasberg, P.; Schaller, G.; Brandes, T.; Esposito, M. Thermodynamics of a physical model implementing a Maxwell demon. *Phys. Rev. Lett.* **2013**, *110*, 040601.
5. Koski, J.V.; Kutvonen, A.; Khaymovich, I.M.; Ala-Nissila, T.; Pekola, J.P. On-Chip Maxwell's Demon as an Information-Powered Refrigerator. *Phys. Rev. Lett.* **2015**, *115*, 260602.
6. Mcgrath, T.; Jones, N.S.; Ten Wolde, P.R.; Ouldridge, T.E. Biochemical Machines for the Interconversion of Mutual Information and Work. *Phys. Rev. Lett.* **2017**, *118*, 028101.
7. Mark, B.L.; Ephraim, Y. An EM algorithm for continuous-time bivariate Markov chains. *Comput. Stat. Data Anal.* **2013**, *57*, 504–517.
8. Ephraim, Y.; Mark, B.L. Bivariate Markov Processes and Their Estimation. *Found. Trends Signal Process.* **2012**, *6*, 1–95.
9. Cover, T.M.; Thomas, J.A. *Elements of Information Theory*, 2nd ed.; John Wiley & Sons: Hoboken, NJ, USA, 2006; ISBN 13-978-0-471-24195-9.
10. Parrondo, J.M.R.; Horowitz, J.M.; Sagawa, T. Thermodynamics of information. *Nat. Phys.* **2015**, *11*, 131–139.
11. Sagawa, T.; Ueda, M. Fluctuation theorem with information exchange: Role of correlations in stochastic thermodynamics. *Phys. Rev. Lett.* **2012**, *109*, 180602.
12. Horowitz, J.M.; Esposito, M. Thermodynamics with Continuous Information Flow. *Phys. Rev. X* **2014**, *4*, 031015.
13. Barato, A.C.; Hartich, D.; Seifert, U. Rate of Mutual Information Between Coarse-Grained Non-Markovian Variables. *J. Stat. Phys.* **2013**, *153*, 460–478.

14.    Wang, J.; Xu, L.; Wang, E.K. Potential landscape and flux framework of nonequilibrium networks: Robustness, dissipation, and coherence of biochemical oscillations. *Proc. Natl. Acad. Sci. USA* **2008**, *105*, 12271–12276.

15.    Wang, J. Landscape and flux theory of non-equilibrium dynamical systems with application to biology. *Adv. Phys.* **2015**, *64*, 1–137.

16.    Li, C.H.; Wang, E.K.; Wang, J. Potential flux landscapes determine the global stability of a Lorenz chaotic attractor under intrinsic fluctuations. *J. Chem. Phys.* **2012**, *136*, 194108.

17.    Schnakenberg, J. Network theory of microscopic and macroscopic behavior of master equation systems. *Rev. Mod. Phys.* **1976**, *48*, 571–585.

18.    Zia, R.K.P.; Schmittmann, B. Probability currents as principal characteristics in the statistical mechanics of non-equilibrium steady states. *J. Stat. Mech.-Theory E* **2007**, *2007*, doi:10.1088/1742-5468/2007/07/p07012.

19.    Maes, C.; Netočný, K. Canonical structure of dynamical fluctuations in mesoscopic nonequilibrium steady states. *Europhys. Lett.* **2008**, *82*, doi:10.1209/0295-5075/82/30003.

20.    Qian, M.P.; Qian, M. Circulation for recurrent markov chains. *Probab. Theory Relat.* **1982**, *59*, 203–210.

21.    Zhang, Z.D.; Wang, J. Curl flux, coherence, and population landscape of molecular systems: Nonequilibrium quantum steady state, energy (charge) transport, and thermodynamics. *J. Chem. Phys.* **2014**, *140*, 245101.

22.    Zhang, Z.D.; Wang, J. Landscape, kinetics, paths and statistics of curl flux, coherence, entanglement and energy transfer in non-equilibrium quantum systems. *New J. Phys.* **2015**, *17*, 043053.

23.    Luo, X.S.; Xu, L.F.; Han, B.; Wang, J. Funneled potential and flux landscapes dictate the stabilities of both the states and the flow: Fission yeast cell cycle. *PLoS Comput. Biol.* **2017**, *13*, e1005710.

24.    Gray, R.; Kieffer, J. Mutual information rate, distortion, and quantization in metric spaces. *IEEE Trans. Inf. Theory* **1980**, *26*, 412–422.

25.    Maes, C.; Redig, F.; van Moffaert, A. On the definition of entropy production, via examples. *J. Math. Phys.* **2000**, *41*, 1528–1554.

26.    Gaspard, P. Time-reversed dynamical entropy and irreversibility in Markovian random processes. *J. Stat. Phys.* **2004**, *117*, 599–615.

27.    Feng, H.D.; Wang, J. Potential and flux decomposition for dynamical systems and non-equilibrium thermodynamics: Curvature, gauge field, and generalized fluctuation-dissipation theorem. *J. Chem. Phys.* **2011**, *135*, 234511.

28.    Polettini, M. Nonequilibrium thermodynamics as a gauge theory. *Europhys. Lett.* **2012**, *97*, 30003.