# Transfer entropy in physical systems and the arrow of time

Richard E. Spinney, Joseph T. Lizier and Mikhail Prokopenko

*Centre for Complex Systems, The University of Sydney, Sydney, New South Wales, Australia, 2006.*

Recent developments have cemented the realization that many concepts and quantities in thermodynamics and information theory are shared. In this paper we consider a highly relevant quantity in information theory and complex systems, the transfer entropy, and explore its thermodynamic role by considering the implications of time reversal upon it. By doing so we highlight the role of information dynamics on the nuanced question of observer perspective within thermodynamics by relating the temporal irreversibility in the information dynamics to the configurational (or spatial) resolution of the thermodynamics. We then highlight its role in perhaps the most enduring paradox in modern physics, the manifestation of a (thermodynamic) arrow of time. We find that for systems that process information such as those undergoing feedback, a robust arrow of time can be formulated by considering both the apparent physical behaviour which leads to conventional entropy production and the information dynamics which leads to a newly defined quantity we call the information theoretic arrow of time. We also offer an interpretation in terms of optimal encoding of observed physical behaviour.

## I. INTRODUCTION

In recent years great progress has been made in describing the thermodynamics of small systems, now increasingly experimentally realizable, through frameworks such as stochastic thermodynamics [1–4]. These frameworks and, more broadly, the constituent work and fluctuation theorems [5–11] have deeply connected entropy production, dissipation, statistical irreversibility and the arrow of time [12–15]. Meanwhile, the field of information theory has enjoyed great success in identifying meaningful components of computation in complex systems strongly implicating transfer entropy amongst other measures [16–22]. More recently still it has become apparent that these two fields, and the concepts they employ, are deeply connected. Through modern consideration of the thermodynamics of systems which were once only explored in the imagination as thought experiments, such as the enduringly compelling Maxwell's demon, information theoretic measures, such as transfer entropy, have taken a central role in *physical* systems [23–35].

Such developments raise intriguing questions: Is information physical? Do the concepts developed studying computation have a role to play in thermodynamics and do thermodynamic concepts have a role in understanding computation? And what of the arrow of time? We observe clouds of smoke billowing from chimneys, not collecting in the atmosphere and flowing into them, plates fall and smash when they hit the floor and the pieces do not spontaneously recombine. This flow of time is considered synonymous with the thermodynamic arrow of time, entropy increases, it does not decrease. Recent advances have formalized the known generalization that such a law, especially for small systems, is only valid statistically [5, 36, 37].

However, even these statistical laws need not hold under feedback [24] or more generally under different observer perspectives of coupled systems. Information theoretic quantities that characterize the information gained by an observer through measurement, or more generally the information that flows out of the the system must be taken into account. These quantities have been used as corrections to the second law, but so far a treatment or generalization of the arrow of time, identified for a system coupled to some feedback mechanism as an involutive quantity that increases, on average, only in the forwards time direction, but decreases by the same amount if viewed in reverse, has not been offered.

In this paper we wish to contribute to these ideas by making further connections between these fields. To do so we seek to understand transfer entropy as completely as possible within physical systems, its role within the question of observer perspective, and to implicate it within a generalized form of the arrow of time. To this end we first establish how to define the transfer entropy in as broad a range of circumstances as regularly arise within non-equilibrium dynamics generalizing it such that it can be defined in the presence of driving and also in the presence of multiple sources of irreversibility. We then, by alluding to the fluctuation theorems, consider the approach of applying time reversal to such a quantity. Using these definitions and a plausible assumption about the dynamics of physically relevant systems we present three main results. The first contrasts the temporal phenomena of irreversibility in the information dynamics with the difference in the implied thermodynamics that results from different configurational resolutions, implicating transfer entropy in the question of observer perspective, uncertainty and irreversibility. This result uniquely considers the difference between transfer entropies in the same direction between partitioned con-

stituents of the system, but with the inclusion of time reversal as opposed to the transfer entropy in different directions as in, for example, [28]. We go on to offer a novel information theoretic interpretation of transfer entropy, relating it to the ability of an observer to encode measured behaviour, identifying it with both a suitable difference in predictive capacity on the target and on the source. Finally, we then bring these concepts together to form a generalized thermodynamic arrow of time and examine the place of such a formalism within the field of information thermodynamics.

## II. TRANSFER ENTROPY IN PHYSICAL SYSTEMS

Introduced by Schreiber [38], the transfer entropy for a pair of coincident time ordered sequences $(\mathbf{x}_0^n, \mathbf{y}_0^n)$ where $\mathbf{x}_0^n = (\mathbf{x}_0, \mathbf{x}_1 \ldots \mathbf{x}_n)$ such that the subscript is a time index is given by

$$T_{\mathbf{y} \to \mathbf{x}}^{(k,l)}(n) = \left\langle \ln \frac{p(\mathbf{x}_n | \mathbf{x}_{n-k}^{n-1}, \mathbf{y}_{n-l}^{n-1})}{p(\mathbf{x}_n | \mathbf{x}_{n-k}^{n-1})} \right\rangle, \quad (1)$$

where $\{k > 0\} \in \mathbb{N}$, $\{l > 0\} \in \mathbb{N}$ and where the bracket notation $\langle \rangle$ indicates averaging over all contained variables including conditional variables by means of path integrals.

Here we generalize such a quantity for use in the physical systems that appear in non-equilibrium thermodynamics formulations such as stochastic thermodynamics [1]. The sequences $(\mathbf{x}_0^n, \mathbf{y}_0^n)$ represent dynamical trajectories of a system which behaves stochastically owing to thermal noise arising from coupling to an idealized environment.

In such formulations one often controls the dynamics through an assumedly deterministic 'protocol' or 'work parameter' [39], $\lambda(t)$, which is an exogenous state variable that parameterizes the system dynamics. By varying the protocol one can drive the system and thus perform or extract work, for example, a piston compressing or expanding a gas. Consequently, the protocol is not a dynamical variable, but a characterization of the probabilities that apply to any transitions and so, for discrete time formulations and approximations, must remain constant during any transition. As such, for a discrete time formulation, the transition $\mathbf{x}_0 \to \mathbf{x}_1$ is parameterized by $\lambda_0$, the transition $\mathbf{x}_1 \to \mathbf{x}_2$ by $\lambda_1$ and so on such that the protocol cannot change using the same time scheme as the system transitions. Typically the protocol parameterizes the dynamics only through its instantaneous value. Despite this, and in the spirit of keeping the treatment of the transfer entropy as a non-parametric statistic, it is instructive to loosen these constraints for initial definitions and apply them again later where appropriate. Consequently the probabilities are generalised $p(\mathbf{x}_n | \mathbf{x}_{n-k}^{n-1}, \mathbf{y}_{n-l}^{n-1}) \to p(\mathbf{x}_n | \mathbf{x}_{n-k}^{n-1}, \mathbf{y}_{n-l}^{n-1}, \lambda_{n-m}^{n-1})$ and

$p(\mathbf{x}_n | \mathbf{x}_{n-k}^{n-1}) \to p(\mathbf{x}_n | \mathbf{x}_{n-k}^{n-1}, \lambda_{n-m}^{n-1})$. However, for brevity, we drop the notational dependence on the protocol (here implicitly assuming that $m = k$) where it is not essential, but all physical transfer entropies should be considered as conditional on the protocol history up to that time. Consequently all probabilities and transition rates should be acknowledged to depend implicitly on the history of the protocol. Further, one can generalise the transfer entropy by allowing for multiple physical mechanisms [40–46], however we defer such a generalization to Appendix A, where it is shown that a failure to account for the additional knowledge of the mechanisms will underestimate the transfer of information.

In this paper we shall, in the main, deal with system in *continuous time* where we write $\mathbf{x}_{t_0}^\tau$ indicating a continuous path function $\{\mathbf{x}(t) : t_0 \leq t < \tau\}$. For these systems it is sensible to consider the *transfer entropy rate* [31, 47, 48]

$$T_{\mathbf{y} \to \mathbf{x}}^{(s,r)}(t) = \lim_{dt \to 0} \frac{1}{dt} \left\langle \ln \frac{p[\mathbf{x}_{t+dt} | \mathbf{x}_{t-s}^t, \mathbf{y}_{t-r}^t]}{p[\mathbf{x}_{t+dt} | \mathbf{x}_{t-s}^t]} \right\rangle \quad (2)$$

where we emphasize the notational scheme where parentheses indicates a function and square brackets indicate a functional. We note the introduction of the positive real valued quantities $s$ and $r$ which, in continuous time, are analogous to the history embedding lengths $k$ and $l$ in discrete time. An underemphasized feature of the transfer entropy in the discussion of thermodynamics is that, just as the entropy production in stochastic thermodynamics has come to be understood as an average of a fluctuating quantity associated within individual realisations, the notion of a *local transfer entropy* exists and is meaningful [17, 20]. Such a quantity behaves analogously to the entropy production for a single realization which is drawn from a distribution which permits negative values. Just as an individual realisation may defy the mean and extract useful energy from its environment, here, knowledge of $\mathbf{y}$ may *misinform* on the level of a single realisation despite the knowledge of $\mathbf{y}$ always leading to improved predictive capabilities on average reflected by a strictly positive mean transfer entropy. Consequently we consider a quantity we call the *pathwise transfer entropy*, $\mathcal{T}_{\mathbf{y} \to \mathbf{x}}^{(s,r)}[\mathbf{x}_{t_0}^\tau, \mathbf{y}_{t_0}^\tau]$, designed to have the property

$$\left\langle \mathcal{T}_{\mathbf{y} \to \mathbf{x}}^{(s,r)}[\mathbf{x}_{t_0}^\tau, \mathbf{y}_{t_0}^\tau] \right\rangle = \int_{t_0}^\tau T_{\mathbf{y} \to \mathbf{x}}^{(s,r)}(t) dt. \quad (3)$$

where $\mathcal{T}_{\mathbf{y} \to \mathbf{x}}^{(s,r)}[\mathbf{x}_{t_0}^\tau, \mathbf{y}_{t_0}^\tau]$, which we write simply as $\mathcal{T}_{\mathbf{y} \to \mathbf{x}}^{(s,r)}$ for brevity, is a fluctuating path dependent functional analogous to the entropy production, heat, work etc. in stochastic thermodynamics. We define such a quantity by considering the log-ratio of two measures over those path realisations.

The first, $\hat{p}[\mathbf{x}_{t_0}^\tau | \mathbf{x}_{t_0-s}^{t_0}, \{\mathbf{y}\}_{t_0-r}^\tau]^{(s,r)}$, designed to replicate the properties of the numerator in eq. (2) at any point

along the path $\mathbf{x}_{t_0}^\tau$, is defined through the limiting procedure

$$\hat{p}[\mathbf{x}_{t_0}^\tau | \mathbf{x}_{t_0-s}^{t_0}, \{\mathbf{y}\}_{t_0-r}^\tau]^{(s,r)} =$$

$$\lim_{\Delta t \to 0} \prod_{i=t_0/\Delta t}^{\tau/\Delta t} p[\mathbf{x}_{(i+1)\Delta t} | \mathbf{x}_{i\Delta t-s}^{i\Delta t}, \mathbf{y}_{i\Delta t-r}^{i\Delta t}]. \quad (4)$$

We call such a quantity the *naive measure*. It is the conditional probability of the path $\mathbf{x}_{t_0}^\tau$ given a *fixed* path $\{\mathbf{y}\}_{t_0-r}^\tau$, where the notation $\{\}$ indicates that the contents is evolved *deterministically*, given historical knowledge of length $s$ and $r$ for the two processes at any point in time along the paths. Alternatively, it can be considered to be the measure that arises when assuming the contents of $\{\}$ as a protocol such that any dependence of $y$ upon $x$ is ignored (i.e. such that $y$ is exogenous to $x$). Whilst there exists no distinction between the naive and actual conditional for a single transition probability, when dealing with measures along complete paths, it is important to note that the naive measure is *not* the conditional measure of the path given $\mathbf{y}_{t_0}^\tau$ [27, 31] i.e.

$$\hat{p}[\mathbf{x}_{t_0}^\tau | \mathbf{x}_{t_0-s}^{t_0}, \{\mathbf{y}\}_{t_0-r}^\tau]^{(s,r)} \neq p[\mathbf{x}_{t_0}^\tau | \mathbf{x}_{t_0-s}^{t_0}, \mathbf{y}_{t_0-r}^\tau]^{(s,r)}. \quad (5)$$

Similarly we construct the measure $p[\mathbf{x}_{t_0}^\tau | \mathbf{x}_{t_0-s}^{t_0}]^{(s)}$, using the same limiting procedure as above, but without dependence on $\mathbf{y}$, which is the conditional probability of the same path without knowledge of $\mathbf{y}$ given historical knowledge of $\mathbf{x}$ of length $s$ at any point along the path. Consequently we have

$$\mathcal{T}_{\mathbf{y} \to \mathbf{x}}^{(s,r)} = \ln \frac{\hat{p}[\mathbf{x}_{t_0}^\tau | \mathbf{x}_{t_0-s}^{t_0}, \{\mathbf{y}\}_{t_0-r}^\tau]^{(s,r)}}{p[\mathbf{x}_{t_0}^\tau | \mathbf{x}_{t_0-s}^{t_0}]^{(s)}}. \quad (6)$$

We note that this functional is, as defined above, evaluated for a path sequence embedded within a larger history such that trajectories before $t = t_0$ can be used in its evaluation. The implication of the above is that a functional of the physical behaviour for the transfer of information can, in principle, be written down and evaluated for individual trajectories. For jump processes in continuous time on a discrete state space such a quantity, for a path starting in configuration $(\mathbf{x}_0, \mathbf{y}_0)$ at time $t_0$ and consisting of $N_x$ transitions, with the $i$th transition being into state $\mathbf{x}_i$ at time $t_i$ is given by

$$\mathcal{T}_{\mathbf{y} \to \mathbf{x}}^{(s,r)} = \sum_{i=1}^{N_x} \ln \frac{\sum_{\mathbf{y}'} W[\mathbf{x}_i, \mathbf{y}' | \mathbf{x}_{t_i-s}^{t_i}, \mathbf{y}_{t_i-r}^{t_i}]}{W[\mathbf{x}_i | \mathbf{x}_{t_i-s}^{t_i}]}$$

$$+ \int_{t_0}^\tau \left( \kappa_x[\mathbf{x}_{t-s}^t] - \kappa_{x|y}[\mathbf{x}_{t-s}^t, \mathbf{y}_{t-r}^t] \right) dt. \quad (7)$$

Here the objects are the joint and coarse grained transition rates

$$W[\mathbf{x}', \mathbf{y}' | \mathbf{x}_{t-s}^t, \mathbf{y}_{t-r}^t] = \lim_{dt \to 0} \frac{1}{dt} p[\mathbf{x}'_{t+dt}, \mathbf{y}'_{t+dt} | \mathbf{x}_{t-s}^t, \mathbf{y}_{t-r}^t]$$

$$W[\mathbf{x}' | \mathbf{x}_{t-s}^t] = \lim_{dt \to 0} \frac{1}{dt} p[\mathbf{x}'_{t+dt} | \mathbf{x}_{t-s}^t]$$

$$= \lim_{dt \to 0} \frac{1}{dt} \frac{\int dy_{t-r}^{t+dt} p[\mathbf{x}_{t-s}^{t+dt}, \mathbf{y}_{t-r}^{t+dt}]}{\int dy_{t-r}^t p[\mathbf{x}_{t-s}^t, \mathbf{y}_{t-r}^t]}, \quad (8)$$

where the notation $\mathbf{x}'_{t+dt}$ indicates the event $\mathbf{x} = \mathbf{x}'$ at time $t + dt$, and the mean escape rates

$$\kappa_x[\mathbf{x}_{t-s}^t] = \sum_{\mathbf{x}' \neq \mathbf{x}_t} W[\mathbf{x}' | \mathbf{x}_{t-s}^t] \quad (9)$$

$$\kappa_{x|y}[\mathbf{x}_{t-s}^t, \mathbf{y}_{t-r}^t] = \sum_{\mathbf{x}' \neq \mathbf{x}_t} \sum_{\mathbf{y}'} W[\mathbf{x}', \mathbf{y}' | \mathbf{x}_{t-s}^t, \mathbf{y}_{t-r}^t]. \quad (10)$$

Derivations for the above can be found in appendix B.

## III. TIME REVERSED TRANSFER ENTROPY

At the heart of modern descriptions of entropy production is time reversal [13]. Within frameworks such as stochastic thermodynamics employing Markovian dynamics a central result, following from the principle of local detailed balance [49–51], is the identification of an external, or medium, entropy production of a single realization as the ratio of conditional measures between the forward path sequence of the realization and the time reversed sequence. We denote the time reversed sequence of a path of duration $\tau$ starting at $t = t_0$, $\mathbf{x}_{t_0}^\tau$, by $\mathbf{x}_{t_0}^{\dagger \tau} = \epsilon \mathbf{x}_\tau^{t_0}$, such that $\mathbf{x}_\tau^{t_0}$ is the sequence $\mathbf{x}_{t_0}^\tau$ 'played' in reverse so that it runs from $\mathbf{x}_\tau$ to $\mathbf{x}_{t_0}$, and where $\epsilon$ is a time reversal operator [52, 53]. For even variables, such as a position, it takes value $\epsilon = +1$ whereas for odd variables, such as a velocity, it takes value $\epsilon = -1$. Similarly the protocol, $\lambda_{t_0}^\tau$, undergoes a time reversal $\lambda_{t_0}^{\dagger \tau} = \epsilon \lambda_\tau^{t_0}$ with odd parity examples being external magnetic fields or torques [54]. For clarity, however, we restrict our discussion to even variables here and throughout. Explicitly, for Markov systems dealing with even variables, in contact with large equilibrium heat baths, we are able to write

$$\frac{p[\mathbf{x}_{t_0}^\tau | \mathbf{x}_{t_0}]}{p[\mathbf{x}_{t_0}^{\dagger \tau} | \mathbf{x}_{t_0}^\dagger]} = \frac{p[\mathbf{x}_{t_0}^\tau | \mathbf{x}_{t_0}]}{p[\mathbf{x}_\tau^{t_0} | \mathbf{x}_\tau]} = \exp\left[ k_B^{-1} \Delta \mathcal{S}_{\text{med}} \right], \quad (11)$$

where we now set $k_B = 1$ throughout. Such an identity powerfully connects thermodynamic entropy with the notion of observed irreversibility and the arrow of time. In models where there is a single environment with a fixed temperature this ratio is readily identified as the heat exported to the bath divided by its temperature. Such a result can be argued more generally [55–59] and allows for a discussion of entropy production in terms of explicit irreversibility and when paired with the principle

of micro-reversibility leads to the celebrated fluctuation theorems [10]. A key feature in the formulation of such results is the inclusion of a system entropy [2] given, for a path of duration $\tau$, starting at time $t = t_0$, by

$$\Delta \mathcal{S}_{\text{sys}} = \ln \frac{p(\mathbf{x}_{t_0})}{p(\mathbf{x}_\tau)} \tag{12}$$

such that the total irreversible entropy production, understood as the total entropy production of the universe consisting of the system and its environment, for a path of duration $\tau$, starting at time $t = t_0$, is given by

$$\begin{aligned}
\Delta \mathcal{S}_{\text{tot}} &= \Delta \mathcal{S}_{\text{sys}} + \Delta \mathcal{S}_{\text{med}} \\
&= \ln \frac{p(\mathbf{x}_{t_0})}{p(\mathbf{x}_\tau)} + \ln \frac{p[\mathbf{x}_{t_0}^\tau | \mathbf{x}_{t_0}]}{p[\mathbf{x}_\tau^{t_0} | \mathbf{x}_\tau]} \\
&= \ln \frac{p[\mathbf{x}_{t_0}^\tau]}{p[\mathbf{x}_\tau^{t_0}]}.
\end{aligned} \tag{13}$$

Here we explore the idea that a meaningful way of connecting information theory with thermodynamics is to apply this notion of time reversal to the quantities central to information theory. To this end we introduce a new quantity, the *time reversed transfer entropy*. We proceed, again, considering only even variables for clarity, by considering the time reversed sequence of a discretized path consisting of $2n$ steps such that it is centred on the $n$th transition. Consequently we have a time reversed path $(\mathbf{x}^\dagger, \mathbf{y}^\dagger)_0^{2n-1} = (\mathbf{x}, \mathbf{y})_{2n-1}^0$, such that $\mathbf{x}_n^\dagger = \mathbf{x}_{n-1}$, and define the *time reversed local transfer entropy* as

$$\begin{aligned}
t_{\mathbf{y} \to \mathbf{x}}^{\dagger, (k,l)}(n) &= \ln \frac{p(\mathbf{x}_n^\dagger | \mathbf{x}_{n-k}^{\dagger \, n-1}, \mathbf{y}_{n-l}^{\dagger \, n-1})}{p(\mathbf{x}_n^\dagger | \mathbf{x}_{n-k}^{\dagger \, n-1})} \\
&= \ln \frac{p(\mathbf{x}_{n-1} | \mathbf{x}_{n+k-1}^n, \mathbf{y}_{n+l-1}^n)}{p(\mathbf{x}_{n-1} | \mathbf{x}_{n+k-1}^n)}.
\end{aligned} \tag{14}$$

The time reversed transfer entropy is then fully defined by the average of this quantity by means of the path probabilities obtained in the forward time direction. In discrete time this is given by

$$\begin{aligned}
T_{\mathbf{y} \to \mathbf{x}}^{\dagger, (k,l)}(n) &= \sum_{\mathbf{x}_{n-1}^{n+k-1}} \sum_{\mathbf{y}_n^{n+l-1}} p(\mathbf{x}_{n-1}^{n+k-1}, \mathbf{y}_n^{n+l-1}) \\
&\quad \times \ln \frac{p(\mathbf{x}_{n-1} | \mathbf{x}_{n+k-1}^n, \mathbf{y}_{n+l-1}^n)}{p(\mathbf{x}_{n-1} | \mathbf{x}_{n+k-1}^n)}
\end{aligned} \tag{15}$$

where the implicit protocol dependence has similarly been time reversed and the mechanisms used match those in the forward path. This is the information transfer understood as a physical quantity, and thus defined using the forward measure, *in the same physical direction*, from $\mathbf{y}$ to $\mathbf{x}$, if time were to run backwards.

Similarly we can build the functional for jump processes in continuous time for the pathwise time reversed

transfer entropy viz.

$$\begin{aligned}
\mathcal{T}_{\mathbf{y} \to \mathbf{x}}^{\dagger, (s,r)} &= \sum_{i=1}^{N_x} \ln \frac{\sum_{\mathbf{y}'} W[\mathbf{x}_{i-1}, \mathbf{y}' | \mathbf{x}_{t_i+s}^{t_i}, \mathbf{y}_{t_i+r}^{t_i}]}{W[\mathbf{x}_{i-1} | \mathbf{x}_{t_i+s}^{t_i}]} \\
&\quad + \int_{t_0}^\tau \left( \kappa_x[\mathbf{x}_{t+s}^t] - \kappa_{x|y}[\mathbf{x}_{t+s}^t, \mathbf{y}_{t+r}^t] \right) dt.
\end{aligned} \tag{16}$$

## IV. RESULT FOR BIPARTITE SYSTEMS

For the remainder of this paper we shall restrict ourselves to classes of system where the entropy production can readily be identified within each component of the system. This restriction is well utilized within the literature [28–30, 42, 60, 61] and amounts to considering Markov bipartite systems. Such a restriction can be summarized by

$$\begin{aligned}
&W[\mathbf{x}', \mathbf{y}' | \mathbf{x}_{t-s}^t, \mathbf{y}_{t-r}^t, \lambda_{t-q}^t] \\
&= \begin{cases}
W(\mathbf{x}', \mathbf{y}_t | \mathbf{x}_t, \mathbf{y}_t, \lambda_t) = W_{\mathbf{x}_t, \mathbf{x}'}^{\mathbf{y}_t}, & \mathbf{x}' \neq \mathbf{x}_t, \mathbf{y}' = \mathbf{y}_t \\
W(\mathbf{x}_t, \mathbf{y}' | \mathbf{x}_t, \mathbf{y}_t, \lambda_t) = W_{\mathbf{x}_t}^{\mathbf{y}_t, \mathbf{y}'}, & \mathbf{x}' = \mathbf{x}_t, \mathbf{y}' \neq \mathbf{y}_t \\
0, & \mathbf{x}' \neq \mathbf{x}_t, \mathbf{y}' \neq \mathbf{y}_t.
\end{cases}
\end{aligned} \tag{17}$$

such that the rates are unchanged in the limit $s \to 0$, $r \to 0$, $q \to 0$ and $W_{\mathbf{x}, \mathbf{x}'}^{\mathbf{y}, \mathbf{y}'} = W(\mathbf{x}', \mathbf{y}' | \mathbf{x}, \mathbf{y}, \lambda)$ is the transition rate from state $(\mathbf{x}, \mathbf{y})$ to $(\mathbf{x}', \mathbf{y}')$ (given $\lambda$) and $W_{\mathbf{x}}^{\mathbf{y}, \mathbf{y}'}$ is shorthand for $W_{\mathbf{x}, \mathbf{x}}^{\mathbf{y}, \mathbf{y}'}$, i.e. a transition where the $x$ state does not change. Explicitly this is the requirement that the joint process is not path dependent beyond the current state and that transitions where both variables change value/state simultaneously are disallowed. One can understand this to be the requirement that both $\mathbf{x}$ and $\mathbf{y}$ can, in principle, be described in isolation such that when they are combined no 'new' transitions are introduced [30]. Alternatively and equivalently one can take this restriction to mean that the thermal noise each experiences is independent from each other. We note that such a structure leads to the important property

$$\begin{aligned}
\sum_{\mathbf{y}'} W(\mathbf{x}', \mathbf{y}' | \mathbf{x}_t, \mathbf{y}_t, \lambda_t) &= W_{\mathbf{x}_t, \mathbf{x}'}^{\mathbf{y}_t} \quad \forall \mathbf{x}' \neq \mathbf{x}_t \\
\sum_{\mathbf{x}'} W(\mathbf{x}', \mathbf{y}' | \mathbf{x}_t, \mathbf{y}_t, \lambda_t) &= W_{\mathbf{x}_t}^{\mathbf{y}_t, \mathbf{y}'} \quad \forall \mathbf{y}' \neq \mathbf{y}_t.
\end{aligned} \tag{18}$$

Since the dynamics are Markov this allows us to evolve the probability density function with the appropriate master equation, once again dropping the notational protocol dependence and noting, for brevity, the contraction $p_{\mathbf{x}}^{\mathbf{y}} = p(\mathbf{x}, \mathbf{y})$,

$$\begin{aligned}
\frac{\partial p_{\mathbf{x}}^{\mathbf{y}}}{\partial t} &= \sum_{\mathbf{x}', \mathbf{y}'} W(\mathbf{x}, \mathbf{y} | \mathbf{x}', \mathbf{y}') p_{\mathbf{x}'}^{\mathbf{y}'} \\
&= \sum_{\mathbf{x}'} W_{\mathbf{x}', \mathbf{x}}^{\mathbf{y}} p_{\mathbf{x}'}^{\mathbf{y}} + \sum_{\mathbf{y}'} W_{\mathbf{x}}^{\mathbf{y}', \mathbf{y}} p_{\mathbf{x}}^{\mathbf{y}'}.
\end{aligned} \tag{19}$$

Identifying the entries

$$W_{\mathbf{x},\mathbf{x}}^{\mathbf{y}} = -\kappa_{x|y}(\mathbf{x},\mathbf{y}) = -\sum_{\mathbf{x}'\neq\mathbf{x}} W_{\mathbf{x},\mathbf{x}'}^{\mathbf{y}}$$

$$W_{\mathbf{x}}^{\mathbf{y},\mathbf{y}} = -\kappa_{y|x}(\mathbf{x},\mathbf{y}) = -\sum_{\mathbf{y}'\neq\mathbf{y}} W_{\mathbf{x}}^{\mathbf{y},\mathbf{y}'}$$

$$W(\mathbf{x},\mathbf{y}|\mathbf{x},\mathbf{y}) = -\kappa_{x,y}(\mathbf{x},\mathbf{y}) = -\kappa_{x|y}(\mathbf{x},\mathbf{y}) - \kappa_{y|x}(\mathbf{x},\mathbf{y}) \tag{20}$$

allows us to view such a equation as continuity equation such that $\dot{p}_{\mathbf{x}}^{\mathbf{y}} = \sum_{\mathbf{x}',\mathbf{y}'} J_{\mathbf{x}'\mathbf{x}}^{\mathbf{y}'\mathbf{y}}$ (with $J_{\mathbf{x}'\mathbf{x}}^{\mathbf{y}'\mathbf{y}} = W(\mathbf{x},\mathbf{y}|\mathbf{x}',\mathbf{y}')p_{\mathbf{x}'}^{\mathbf{y}'}$) being a sum over local probability currents [30]. The key property of the system being bipartite is that all other parts of the system are stationary whilst another part transitions. Consequently, as the second line of eq. (19) demonstrates, in these systems, the total probability current separates into distinct components. As a consequence it is well known that the mean rate of total entropy production of the universe (the system and its environment) separates [28–30] such that

$$\dot{S}_{\text{tot}} = \frac{d\langle\Delta\mathcal{S}_{\text{tot}}\rangle}{dt} = \frac{d\langle\Delta\mathcal{S}_{\text{sys}}\rangle}{dt} + \frac{d\langle\Delta\mathcal{S}_{\text{med}}\rangle}{dt}$$

$$= \sum_{\mathbf{y},\mathbf{x}}\sum_{\mathbf{x}'\neq\mathbf{x},\mathbf{y}'\neq\mathbf{y}} p_{\mathbf{x}}^{\mathbf{y}} W_{\mathbf{x},\mathbf{x}'}^{\mathbf{y},\mathbf{y}'}\left(\ln\frac{p_{\mathbf{x}}^{\mathbf{y}}}{p_{\mathbf{x}'}^{\mathbf{y}'}} + \ln\frac{W_{\mathbf{x},\mathbf{x}'}^{\mathbf{y},\mathbf{y}'}}{W_{\mathbf{x}',\mathbf{x}}^{\mathbf{y}',\mathbf{y}}}\right)$$

$$= \sum_{\mathbf{y},\mathbf{x}}\sum_{\mathbf{x}'\neq\mathbf{x}} p_{\mathbf{x}}^{\mathbf{y}} W_{\mathbf{x},\mathbf{x}'}^{\mathbf{y}}\ln\frac{p_{\mathbf{x}}^{\mathbf{y}} W_{\mathbf{x},\mathbf{x}'}^{\mathbf{y}}}{p_{\mathbf{x}'}^{\mathbf{y}} W_{\mathbf{x}',\mathbf{x}}^{\mathbf{y}}}$$

$$+ \sum_{\mathbf{x},\mathbf{y}}\sum_{\mathbf{y}'\neq\mathbf{y}} p_{\mathbf{x}}^{\mathbf{y}} W_{\mathbf{x}}^{\mathbf{y},\mathbf{y}'}\ln\frac{p_{\mathbf{x}}^{\mathbf{y}} W_{\mathbf{x}}^{\mathbf{y},\mathbf{y}'}}{p_{\mathbf{x}}^{\mathbf{y}'} W_{\mathbf{x}}^{\mathbf{y}',\mathbf{y}}}$$

$$= \frac{d\langle\Delta\mathcal{S}_i^x\rangle}{dt} + \frac{d\langle\Delta\mathcal{S}_i^y\rangle}{dt} = \dot{S}_i^x + \dot{S}_i^y \tag{21}$$

with each of $\dot{S}_i^x$ and $\dot{S}_i^y$ bounded from below by 0 by the log sum inequality. This ability to separate the probability current can be thought of as arising from the ability to treat the other parts of the system, $\mathbf{y}$, as a protocol in that they entirely parameterize the dynamics experienced by $\mathbf{x}$ throughout a transition. This allows us to untangle the joint dynamics and thermodynamic quantities like the exported entropy into contributions entirely attributable to transitions in either $\mathbf{x}$ or $\mathbf{y}$. For example the exported entropy production for a realization is found by constructing eq. (11).

By using the form of conditional path probabilities found in eq. (B4) of Appendix B, whilst recognizing eq. (18) and the Markov property

$$\kappa_{x|y}[\mathbf{x}_{t-s}^t, \mathbf{y}_{t-r}^t] = \kappa_{x|y}[\mathbf{x}_{t+s}^t, \mathbf{y}_{t+r}^t] = \kappa_{x|y}(\mathbf{x}_t, \mathbf{y}_t), \tag{22}$$

such that the mean escape rates are equivalent in both the forward and reverse time directions one finds, for

$N = N_x + N_y$ transitions into states $(\mathbf{x}_k, \mathbf{y}_k)$ at times $t_k$

$$\Delta\mathcal{S}_{\text{med}} = \sum_{k=1}^{N}\ln\frac{W(\mathbf{x}_k,\mathbf{y}_k|\mathbf{x}_{k-1},\mathbf{y}_{k-1},\lambda_{t_k})}{W(\mathbf{x}_{k-1},\mathbf{y}_{k-1}|\mathbf{x}_k,\mathbf{y}_k,\lambda_{t_k})}$$

$$= \sum_{i=1}^{N_x}\ln\frac{W_{\mathbf{x}_{i-1},\mathbf{x}_i}^{\mathbf{y}_{i-1}}}{W_{\mathbf{x}_i,\mathbf{x}_{i-1}}^{\mathbf{y}_{i-1}}} + \sum_{j=1}^{N_y}\ln\frac{W_{\mathbf{x}_{j-1}}^{\mathbf{y}_{j-1},\mathbf{y}_j}}{W_{\mathbf{x}_{j-1}}^{\mathbf{y}_j,\mathbf{y}_{j-1}}}$$

$$= \Delta\mathcal{S}_{\text{med}}^x + \Delta\mathcal{S}_{\text{med}}^y. \tag{23}$$

Here each contribution is formed with a separate sum using its own time scheme which only counts transitions in that subsystem. Explicitly, there are $N_x$ transitions into states $\mathbf{x}_i$ occurring at times $t_i$, with $\mathbf{x}_{i-1}$ and $\mathbf{y}_{i-1}$ being the immediately preceding values and $N_y$ transitions into states $\mathbf{y}_j$ occurring at times $t_j$ again with $\mathbf{x}_{j-1}$ and $\mathbf{y}_{j-1}$ being the immediately preceding values.

In contrast to the above, when considering the coarse grained dynamics, for instance in the absence of knowledge of $\mathbf{y}$, it is important to appreciate that the behaviour of $\mathbf{x}$ is still, in general, non-Markov in that inferences can be made using its path history as well as historical values of the protocol. Consequently, for the coarse grained system, there is no general simplification of the transition rates and we retain the form $W[\mathbf{x}'|\mathbf{x}_{t-s}^t, \lambda_{t-q}^t]$, but drop the notational dependence on the protocol assuming $q = s$.

By using the above properties, particularly eqs. (23) and (18), we arrive at our first result for Markovian bipartite systems in continuous time by subtracting eq. (16) from eq. (7)

$$\mathcal{T}_{\mathbf{y}\to\mathbf{x}}^{(s)} - \mathcal{T}_{\mathbf{y}\to\mathbf{x}}^{\dagger,(s)} = \sum_{i=1}^{N_x}\ln\frac{W_{\mathbf{x}_{i-1},\mathbf{x}_i}^{\mathbf{y}_{i-1}}}{W_{\mathbf{x}_i,\mathbf{x}_{i-1}}^{\mathbf{y}_{i-1}}} - \sum_{i=1}^{N_x}\ln\frac{W[\mathbf{x}_i|\mathbf{x}_{t_i-s}^{t_i}]}{W[\mathbf{x}_{i-1}|\mathbf{x}_{t_i+s}^{t_i}]}$$

$$- \int_0^\tau\left(\kappa_x[\mathbf{x}_{t+s}^t] - \kappa_x[\mathbf{x}_{t-s}^t]\right)dt$$

$$= \Delta\mathcal{S}_{\text{med}}^x - \ln\frac{p[\mathbf{x}_{t_0}^\tau|\mathbf{x}_{t_0-s}^{t_0}]^{(s)}}{p[\mathbf{x}_\tau^{t_0}|\mathbf{x}_{\tau+s}^\tau]^{(s)}}$$

$$\mathcal{T}_{\mathbf{y}\to\mathbf{x}}^{(s)} - \mathcal{T}_{\mathbf{y}\to\mathbf{x}}^{\dagger,(s)} = \Delta\mathcal{S}_{\text{med}}^x - \Delta\mathbb{S}_{\text{med}}^{x,(s)}. \tag{24}$$

We note that the parameter $r$ has no effect and has thus been dropped since the joint process is assumed to be Markov. Here the first term is identified as the $x$ component of eq. (23) and the remainder is the ratio of unconditioned path probabilities of generic continuous time processes in $\mathbf{x}$, in the forward and reverse time directions given by eq. (B4) in Appendix B, given historical knowledge $s$ at all points along the paths. We emphasize the non-Markovian nature of the probabilities appearing in the final term and point out that, as with eq. (6), the above quantity is formulated for a sequence embedded in a wider historical (and future) sequence thus leading to probabilities conditional on previous and future path sequences.

On the other hand we may consider cases where there is a definite time origin and horizon such that no information about the path is available before $t_0$ or after $\tau$. In these cases the path measures reduce to the form $p[\mathbf{x}_{t_0}^\tau|\mathbf{x}_{t_0}]^{(s)}$ where at each point $t \in [t_0, \tau]$, $\min(t - t_0, s)$ seconds of path history are utilised in its formulation. Of particular note are such situations, but where we otherwise let $s \to \infty$ such that at any point, $t$, we may consider $s$ taking the value $t - t_0$. Similarly in the reverse measures we may consider $s$ taking the value $\tau - t$. We write the pathwise and reverse pathwise transfer entropy under these conditions $\mathcal{T}_{\mathbf{y}\to\mathbf{x}}^{(t-t_0)} = \mathcal{T}_{\mathbf{y}\to\mathbf{x}}$ and $\mathcal{T}_{\mathbf{y}\to\mathbf{x}}^{\dagger,(\tau-t)} = \mathcal{T}_{\mathbf{y}\to\mathbf{x}}^\dagger$. These quantities are of particular importance because they use the natural measures $p[\mathbf{x}_{t_0}^\tau|\mathbf{x}_{t_0}]$ that result from marginalizing the joint measure $p[\mathbf{x}_{t_0}^\tau, \mathbf{y}_{t_0}^\tau|\mathbf{x}_{t_0}, \mathbf{y}_{t_0}]$ and as such appear in our subsequent results. Under these circumstances we are free to include the coarse grained system entropy in $\mathbf{x}$, using the form appearing in eq. (12), but with distributions over $\mathbf{x}$ arising from the marginalized joint distributions over $\mathbf{x}$ and $\mathbf{y}$, to both entropy productions thus arriving at

$$\mathcal{T}_{\mathbf{y}\to\mathbf{x}} - \mathcal{T}_{\mathbf{y}\to\mathbf{x}}^\dagger = \Delta\mathcal{S}_{\text{tot}}^x - \Delta\mathbb{S}_{\text{tot}}^x \qquad (25)$$

where $\Delta\mathcal{S}_{\text{tot}}^x = \Delta\mathbb{S}_{\text{sys}}^x + \Delta\mathcal{S}_{\text{med}}^x$ and where $\Delta\mathbb{S}_{\text{sys}}^x = \ln(p(\mathbf{x}_{t_0})/p(\mathbf{x}_\tau))$. Here and throughout the notation $\mathbb{S}$ indicates that the quantity is formulated by considering coarse grained distributions and path probabilities.

Eqs. (24) and (25) amount to a precise relation between *information dynamics*, in the same physical direction, in both the forward and reverse time directions and the *irreversibility* of the subsystem $\mathbf{x}$ when observed at different levels of description or, with an appropriate partition [44], length scales. The last term in eq. (25) is a quantity known as the *coarse grained entropy production* [28, 62, 63] which denotes the irreversibility of the subsystem $\mathbf{x}$ when no information about $\mathbf{y}$, not even necessarily its existence, is known. When the process $\mathbf{y}$ is much faster than $\mathbf{x}$ then the coarse grained entropy production becomes a Markovian entropy production [41]. In contrast $\Delta\mathcal{S}_{\text{tot}}^x$ can be interpreted as the entropy production of $\mathbf{x}$ when $\mathbf{y}$ is known to exist, but the details of which are not known [30], or as the entropy production of $\mathbf{x}$ under the assumption $\mathbf{y}$ is a protocol. It is thus worth pointing out that there are therefore three distinct entropy productions, in the form of system entropies represented by logarithmic one time probabilities and apparent medium entropies represented by conditional path probabilities, relevant to distinct observer perspectives when viewing the system

$$\Delta\mathcal{S}_{\text{tot}} = \Delta\mathcal{S}_{\text{sys}} + \Delta\mathcal{S}_{\text{med}} = \ln\frac{p[\mathbf{x}_{t_0}^\tau, \mathbf{y}_{t_0}^\tau]}{p[\mathbf{x}_\tau^{t_0}, \mathbf{y}_\tau^{t_0}]}$$

$$\Delta\mathcal{S}_{\text{tot}}^x = \Delta\mathbb{S}_{\text{sys}}^x + \Delta\mathcal{S}_{\text{med}}^x = \ln\frac{p(\mathbf{x}_{t_0})\hat{p}[\mathbf{x}_{t_0}^\tau|\mathbf{x}_{t_0}, \{\mathbf{y}\}_{t_0}^\tau]}{p(\mathbf{x}_\tau)\hat{p}[\mathbf{x}_\tau^{t_0}|\mathbf{x}_\tau, \{\mathbf{y}\}_\tau^{t_0}]}$$

$$= \ln\frac{\hat{p}[\mathbf{x}_{t_0}^\tau|\{\mathbf{y}\}_{t_0}^\tau]}{\hat{p}[\mathbf{x}_\tau^{t_0}|\{\mathbf{y}\}_\tau^{t_0}]}$$

$$\Delta\mathbb{S}_{\text{tot}}^x = \Delta\mathbb{S}_{\text{sys}}^x + \Delta\mathbb{S}_{\text{med}}^x = \ln\frac{p[\mathbf{x}_{t_0}^\tau]}{p[\mathbf{x}_\tau^{t_0}]}. \qquad (26)$$

The first represents the irreversibility in the total system, the second, based on the naive measure, is the apparent irreversibility when focused on $\mathbf{x}$ and treating $\mathbf{y}$ as a protocol and the third the irreversibility of $\mathbf{x}$ when $\mathbf{y}$ is not known to or measured by the observer. We emphasize that $\Delta\mathcal{S}_{\text{tot}}^x \neq \Delta\mathcal{S}_i^x$, differing, in the mean, by a quantity called the information flow [30, 31] defined as the mutual information rate arising from the probability current in $\mathbf{x}$. Alternatively, and equivalently, in a steady state, this difference, in addition to the coarse grained system entropy rate of $\mathbf{x}$ has been identified as the learning rate or entropy reduction [64].

The result in eq. (25) can also be seen as a precise division, into physical directions, of information transfers, for finite processes, of the so called *mutual entropy production* appearing in [28]. This relies on the property that for bipartite systems the trajectory mutual information rate comprises the transfer entropy rate in opposite directions [29, 31]. Explicitly, using the above formalism it is equivalent to the result

$$\mathcal{T}_{\mathbf{y}\to\mathbf{x}} - \mathcal{T}_{\mathbf{y}\to\mathbf{x}}^\dagger + \mathcal{T}_{\mathbf{x}\to\mathbf{y}} - \mathcal{T}_{\mathbf{x}\to\mathbf{y}}^\dagger =$$
$$- I_m(\mathbf{x}_{t_0}, \mathbf{y}_{t_0}) + \Delta\mathcal{S}_{\text{tot}} - \Delta\mathbb{S}_{\text{tot}}^x - \Delta\mathbb{S}_{\text{tot}}^y \qquad (27)$$

where $I_m(\mathbf{x}_{t_0}, \mathbf{y}_{t_0}) = \ln p(\mathbf{x}_{t_0}, \mathbf{y}_{t_0})/p(\mathbf{x}_{t_0})p(\mathbf{y}_{t_0})$ is the initial local mutual information between $\mathbf{x}$ and $\mathbf{y}$. The result in [28] then concerns the mean transfer entropy and entropy production *rates* such that the contribution from the initial mutual information becomes negligible.

### A. Information transfer as an encoding cost

We now consider an alternative definition for the transfer entropy in physical bipartite systems. Considering again the natural measures given a fixed time origin at $t = t_0$ such that $s = t - t_0$ for all times, we now point out that for Markovian bipartite systems we have the property

$$p[\mathbf{x}_{t_0}^\tau, \mathbf{y}_{t_0}^\tau] = p(\mathbf{x}_{t_0}, \mathbf{y}_{t_0})\hat{p}[\mathbf{x}_{t_0}^\tau|\mathbf{x}_{t_0}, \{\mathbf{y}\}_{t_0}^\tau]\hat{p}[\mathbf{y}_{t_0}^\tau|\mathbf{y}_{t_0}, \{\mathbf{x}\}_{t_0}^\tau]$$

$$= \frac{p(\mathbf{x}_{t_0}, \mathbf{y}_{t_0})}{p(\mathbf{x}_{t_0})p(\mathbf{y}_{t_0})}\hat{p}[\mathbf{x}_{t_0}^\tau|\{\mathbf{y}\}_{t_0}^\tau]\hat{p}[\mathbf{y}_{t_0}^\tau|\{\mathbf{x}\}_{t_0}^\tau]$$

$$= p[\mathbf{y}_{t_0}^\tau|\mathbf{x}_{t_0}^\tau]p[\mathbf{x}_{t_0}^\tau]. \qquad (28)$$

The first line underlies eq. (27) and is recognizable by considering the general form of conditional path probabilities given by eq. (B4) of Appendix B along with the properties in eq. (18) and the divisibility of the total mean escape rate in eq. (20). Consequently we may write

$$\mathcal{T}_{\mathbf{y}\to\mathbf{x}} = \ln\frac{\hat{p}[\mathbf{x}_{t_0}^\tau|\{\mathbf{y}\}_{t_0}^\tau]}{p[\mathbf{x}_{t_0}^\tau]}$$
$$= \ln\frac{p[\mathbf{y}_{t_0}^\tau|\mathbf{x}_{t_0}^\tau]}{\hat{p}[\mathbf{y}_{t_0}^\tau|\{\mathbf{x}\}_{t_0}^\tau]} - I_m(\mathbf{x}_{t_0},\mathbf{y}_{t_0}) \qquad (29)$$

where $I_m(\mathbf{x}_{t_0},\mathbf{y}_{t_0})$ is the local mutual information at the start of the process. This relates the transfer entropy, usually described in terms of different probability measures on the target, in terms of probability measures on the source. Specifically we interpret the remaining term in eq. (29) as a pointwise excess, that is additional, *code length* [65] that is associated with encoding the behaviour of subsystem $\mathbf{y}$ given knowledge of the behaviour of subsystem $\mathbf{x}$ under the naive measure, the assumption that $\mathbf{x}$ is deterministic or that is $\mathbf{x}$ is a classical switching protocol. Explicitly, if we encode messages, optimally, based on the assumption that they arrive with probabilities $p_e$, but they actually arrive according to different probabilities $p_a$ the mean additional cost of encoding messages given this error, in nats, is $D_{KL}(p_a||p_e) = \sum p_a \ln(p_a/p_e)$. The pointwise cost, $\ln(p_a/p_e)$, associated with a single message is the quantity we discuss here. Such an interpretation has been considered before, again with implications for the arrow of time, in systems without feedback [66]. We emphasize that the encoding task, and thus additional code length from a non-optimal encoding, is to encode the dynamical sequence $\mathbf{y}_{t_0}^\tau$ given the knowledge that the realization in $\mathbf{x}$ is $\mathbf{x}_{t_0}^\tau$ and that $\mathbf{x}$ and $\mathbf{y}$ are not to be interpreted as channel sources and outputs as is often considered in the classical information theory literature [65]. We also mention that discrete systems can be losslessly encoded by representing the paths using tuples of the state variable and the times between transitions, whereas for fully continuous systems the encoding is deemed to be approximate being achieved by some representation, i.e. [67], given an available bit rate which can be considered in the limit where relationships between code lengths hold. This allows us to write our second main result

$$\mathcal{T}_{\mathbf{y}\to\mathbf{x}} = d(p_{\mathbf{y}|\mathbf{x}}||\hat{p}_{\mathbf{y}|\{\mathbf{x}\}}) - I_m(\mathbf{x}_{t_0},\mathbf{y}_{t_0}). \qquad (30)$$

Here $d(p_{\mathbf{y}|\mathbf{x}}||\hat{p}_{\mathbf{y}|\{\mathbf{x}\}}) = \ln p[\mathbf{y}_{t_0}^\tau|\mathbf{x}_{t_0}^\tau]/\hat{p}[\mathbf{y}_{t_0}^\tau|\{\mathbf{x}\}_{t_0}^\tau]$ is the additional code length with notation chosen to reflect its form as a pointwise (i.e. local) Kullback-Leibler divergence. We see that the transfer entropy and additional code length are exactly equal given an initial uncorrelated state. We also point out that the *transfer entropy rate*, relevant in systems converged upon a steady state, defined in the limit $t \to \infty$, such that the initial mutual information becomes negligible, can now be defined

through two equivalent statements

$$T_{y\to x} = \lim_{\tau\to\infty}\frac{1}{\tau}\left\langle \ln\frac{\hat{p}[\mathbf{x}_{t_0}^\tau|\{\mathbf{y}\}_{t_0}^\tau]}{p[\mathbf{x}_{t_0}^\tau]}\right\rangle$$
$$= \lim_{\tau\to\infty}\frac{1}{\tau}\left\langle \ln\frac{p[\mathbf{y}_{t_0}^\tau|\mathbf{x}_{t_0}^\tau]}{\hat{p}[\mathbf{y}_{t_0}^\tau|\{\mathbf{x}\}_{t_0}^\tau]}\right\rangle \qquad (31)$$

thus identifying the transfer entropy rate as *both* the increase in predictability in the *target* given knowledge of the source, as is the commonly understood definition, but also as the additional cost associated with encoding the *source* given knowledge of the target when treating the target as if it were an exogenous variable i.e. a protocol.

### B. Discrete time systems

It should be noted that the above result requires the system to be bipartite, Markov and in continuous time. In continuous time a lack of correlation in the thermal noise experienced by the two subsystems and an ability to divide the current into two components coincide. However this is not the case in discrete time and so the result is required to be adapted in such cases depending on which property is chosen. One can assert that the transition probabilities experienced by subsystem $\mathbf{x}$ should not depend on the future state of subsystem $\mathbf{y}$. We call such systems separable and have the property

$$p(\mathbf{x}_i,\mathbf{y}_i|\mathbf{x}_{i-1},\mathbf{y}_{i-1}) = p(\mathbf{x}_i|\mathbf{x}_{i-1},\mathbf{y}_{i-1})p(\mathbf{y}_i|\mathbf{x}_{i-1},\mathbf{y}_{i-1}). \qquad (32)$$

Such a property leads to a medium entropy production of the form

$$\Delta\mathcal{S}_{\text{med}} = \ln\frac{p(\mathbf{x}_i|\mathbf{x}_{i-1},\mathbf{y}_{i-1})}{p(\mathbf{x}_{i-1}|\mathbf{x}_i,\mathbf{y}_i)} + \ln\frac{p(\mathbf{y}_i|\mathbf{x}_{i-1},\mathbf{y}_{i-1})}{p(\mathbf{y}_{i-1}|\mathbf{x}_i,\mathbf{y}_i)}$$
$$= \Delta\tilde{\mathcal{S}}_{\text{med}}^x + \Delta\tilde{\mathcal{S}}_{\text{med}}^y \qquad (33)$$

where $\Delta\tilde{\mathcal{S}}_{\text{med}}^x$ is a component of the exported entropy production associated with a transition in $\mathbf{x}$, but cannot be associated exclusively with the current in $\mathbf{x}$. For such systems eq. (24) becomes

$$\mathcal{T}_{\mathbf{y}\to\mathbf{x}}^{(s)} - \mathcal{T}_{\mathbf{y}\to\mathbf{x}}^{\dagger,(s)} = \Delta\tilde{\mathcal{S}}_{\text{med}}^x - \Delta\mathbb{S}_{\text{med}}^{x,(s)} \qquad (34)$$

with $\Delta\tilde{\mathcal{S}}_{\text{med}}^x$ coinciding with $\Delta\mathcal{S}_{\text{med}}^x$ if a time step $\Delta t$ is associated with each transition and that time step is taken to 0.

Alternatively one can assert that the system is bipartite similarly to the systems which eq. (24) applies to in the sense that both $\mathbf{x}$ and $\mathbf{y}$ cannot change state within one time step. In such systems all quantities can be divided into two components associated with $\mathbf{x}$ and $\mathbf{y}$ currents and transitions [30]. To select the component that is associated with the $\mathbf{x}$ current component of a quantity we introduce the operator $\hat{\mathcal{X}}$ which filters out all steps

that involve $\mathbf{y}$ transitions such that it applies the delta function $\delta_{\mathbf{y}_i \mathbf{y}_{i+1}}$ at each time step. In such systems eq. (24) becomes

$$\hat{\mathcal{X}}\left[\mathcal{T}_{\mathbf{y}\to\mathbf{x}}^{(s)} - \mathcal{T}_{\mathbf{y}\to\mathbf{x}}^{\dagger,(s)}\right] = \Delta\mathcal{S}_{\text{med}}^x - \Delta\mathbb{S}_{\text{med}}^{x,(s)}$$
$$= \hat{\mathcal{X}}\left[\Delta\mathcal{S}_{\text{med}} - \Delta\mathbb{S}_{\text{med}}^{x,(s)}\right]. \quad (35)$$

## V. GENERALIZING THE ARROW OF TIME FOR CONTINUOUS FEEDBACK PROCESSES

It is well known that in the absence of feedback or coupling to hidden variables the statement of the second law is a mathematical encoding of the apparent time reversal asymmetry that we observe in physical phenomena. Indeed, a famous and important corollary is that in equilibrium the absence of entropy production is to be directly associated with an inability to perform any test that can determine whether time is running forwards or backwards [39]. Indeed questions around how one might be able to make such a decision given finite observations are both interesting and contemporary [68, 69].

However, with feedback, or more generally, with focus only on a portion of the system, an observer might be systematically misinformed in their conclusion about which direction time is flowing. In such a set up one can take the view that an observer would conflate the motion of the rest of the system with a protocol which, assumedly, but erroneously, did not depend on the behaviour of the subsystem. Alternatively one can argue that the observer only has the ability to measure the subsystem in question, but can measure the entropy exported by that subsystem which implicitly depends on the total system's state. The entropy production with such incomplete knowledge would, for subsystem $\mathbf{y}$, be given by [30]

$$\Delta\mathcal{S}_{\text{tot}}^y = \Delta\mathbb{S}_{\text{sys}}^y + \Delta\mathcal{S}_{\text{med}}^y$$
$$= \ln\frac{p(\mathbf{y}_{t_0})}{p(\mathbf{y}_\tau)} + \sum_{j=1}^{N_y}\ln\frac{W_{\mathbf{x}_{j-1}}^{\mathbf{y}_{j-1},\mathbf{y}_j}}{W_{\mathbf{x}_{j-1}}^{\mathbf{y}_j,\mathbf{y}_{j-1}}} = \ln\frac{\hat{p}[\mathbf{y}_{t_0}^\tau|\{\mathbf{x}\}_{t_0}^\tau]}{\hat{p}[\mathbf{y}_\tau^{t_0}|\{\mathbf{x}\}_\tau^{t_0}]} \quad (36)$$

which has no bounds on its sign since the average

$$\langle\Delta\mathcal{S}_{\text{tot}}^y\rangle = \int d\mathbf{x}_{t_0}^\tau p[\mathbf{x}_{t_0}^\tau]\int d\mathbf{y}_{t_0}^\tau p[\mathbf{y}_{t_0}^\tau|\mathbf{x}_{t_0}^\tau]\ln\frac{\hat{p}[\mathbf{y}_{t_0}^\tau|\{\mathbf{x}\}_{t_0}^\tau]}{\hat{p}[\mathbf{y}_\tau^{t_0}|\{\mathbf{x}\}_\tau^{t_0}]} \quad (37)$$

does not obey the log sum inequality because $p[\mathbf{y}_{t_0}^\tau|\mathbf{x}_{t_0}^\tau] \neq \hat{p}[\mathbf{y}_{t_0}^\tau|\{\mathbf{x}\}_{t_0}^\tau]$ owing to the general property that the evolution of $\mathbf{x}$ also depends on $\mathbf{y}$ [27]. This gives insight into the apparent arrow of time in such systems: the entropy production encodes the preference for the forward process over the corresponding reverse process when the motion $\mathbf{x}$ is assumed to be

a deterministic protocol such that is predetermined and independent of $\mathbf{y}$; an assumption facilitated by the bipartite structure of the dynamics.

In language more fitting with our everyday experience of the arrow of time, based on our pre-existing knowledge of the behaviour of matter (the *known*, albeit perhaps imperfectly, probabilities of paths given fixed protocols), we are able to discern which way a film is being played based on our level of 'surprise' in each direction: removing a barrier (a deterministic protocol) that releases coloured dye into a fluid is deemed probable whereas the recollecting of dye immediately before the closure of the barrier (the reversed protocol) is deemed improbable. But such an observation for a small enough system with enough thermal noise and with a sensing protocol might be engineered to occur more frequently. Such realizations are associated with negative entropy productions and with feedback the entropy might be negative on average. Using this observed entropy production could lead to incorrect inferences about the arrow of time: on average the time reversed behaviour could be more probable than the forward behaviour.

There are now, well known and appealing, results concerning feedback that deal with the sum of the exported entropy and exported information. The exported information in these results is measured through an information theoretic quantity that generalizes to the pathwise transfer entropy in the continuous feedback limit, which must be greater than zero [23, 27]. However, there is no direct analogy to the canonical case of time asymmetry whereby a single involutive quantity, measured in each time direction, either increases or decreases in the mean [39, 66, 70] as such results require direct external interference in the system with the requirement that no feedback occurs in the reverse time direction. To proceed we imagine a different approach and recognize that, heuristically, the arrow of time is based on assigning a physically observed weight (probability) to the forward and reverse behaviour. We propose that we augment the arrow of time by including an information theoretic weight to each of these behaviours. To do so we introduce a new quantity designed to be an information theoretic analogue to the usual definition of the entropy production appearing in eq. (36)

$$\Delta\mathcal{S}_{\text{info}}^{y|x} = \ln\frac{\frac{p(\mathbf{x}_{t_0}|\mathbf{y}_{t_0})}{p(\mathbf{x}_{t_0})}}{\frac{p(\mathbf{x}_\tau|\mathbf{y}_\tau)}{p(\mathbf{x}_\tau)}} + \ln\frac{\frac{\hat{p}[\mathbf{x}_{t_0}^\tau|\mathbf{x}_{t_0},\{\mathbf{y}\}_{t_0}^\tau]}{p[\mathbf{x}_{t_0}^\tau|\mathbf{x}_{t_0}]}}{\frac{\hat{p}[\mathbf{x}_\tau^{t_0}|\mathbf{x}_\tau,\{\mathbf{y}\}_\tau^{t_0}]}{p[\mathbf{x}_\tau^{t_0}|\mathbf{x}_\tau]}}$$
$$= \ln\frac{\exp\left[I_m(\mathbf{x}_{t_0},\mathbf{y}_{t_0})\right]}{\exp\left[I_m(\mathbf{x}_\tau,\mathbf{y}_\tau)\right]} + \ln\frac{\exp\left[\mathcal{T}_{\mathbf{y}\to\mathbf{x}}\right]}{\exp\left[\mathcal{T}_{\mathbf{y}\to\mathbf{x}}^\dagger\right]}$$
$$= I_m(\mathbf{x}_{t_0},\mathbf{y}_{t_0}) + \mathcal{T}_{\mathbf{y}\to\mathbf{x}} - I_m(\mathbf{x}_\tau,\mathbf{y}_\tau) - \mathcal{T}_{\mathbf{y}\to\mathbf{x}}^\dagger. \quad (38)$$

Here $I_m$ is the local mutual information and plays the role of the system entropy in eq. (36) which conditions

the path probabilities for the entropy production, but here characterizes the initial state of the shared information. It should be noted that in previous approaches, the transfer entropy is arrived at by comparing probabilities from a forward experiment and a carefully chosen reverse experiment which implements no feedback (see eqs. (55), (63) and (65) in [27]). As we will see, such a situation can be equivalently characterized by a vanishing *reverse* transfer entropy. In contrast here, the information transfer that occurs autonomously, $\mathcal{T}_{\mathbf{y}\to\mathbf{x}}$, is considered directly. This is then contrasted with the autonomous information transfer in the reverse time direction, $\mathcal{T}^{\dagger}_{\mathbf{y}\to\mathbf{x}}$. It should also be noted that $\Delta\mathcal{S}^{y|x}_{\mathrm{info}}$ is constructed so that there is an explicit origin and horizon for the paths $(\mathbf{x}^{\tau}_{t_0}, \mathbf{y}^{\tau}_{t_0})$ such that the local transfer entropies cannot utilize any path histories before $t = t_0$ or beyond $t = \tau$ and with $s$ set to $t - t_0, \tau - t$ for the forward and reverse processes such that all available information is utilized.

We interpret eq. (38) as characterizing the difference between the initial shared information plus the information that is exported by $y$ (though, unlike the medium entropy production, may be possibly still dynamically shared with $y$) and the final shared information plus the information that is exported by $y$ in the reverse process. We argue that this quantity effects an information theoretic arrow of time: it contrasts the *information dynamics* in the forward and reverse time directions whilst, eq. (36), which we consider to be the physical arrow of time contrasts the apparent *physical dynamics* in the forward and reverse time directions. Eq. (38) is similar to the quantity introduced in [71], but includes the newly introduced reverse transfer entropy. Similarly to the observed entropy production in eq. (36), in the steady state the boundary terms do not contribute on average so $\langle\Delta\mathcal{S}^{y|x}_{\mathrm{info}}\rangle$ becomes a balance between the information transferred in each time direction which is to be contrasted with $\langle\Delta\mathcal{S}^{y}_{\mathrm{tot}}\rangle$ which is a balance between the apparent path likelihoods in each time direction. Crucial to this concept is the idea that information transfer can be understood physically, that it is formed from pre-existing known physical probability measures. This means that, just as with the physical arrow of time, probabilities are not determined empirically in the direction the film happens to be playing (a technique which would leave us unable to detect which way time was flowing with or without feedback), but based on established physical behaviour when time is known to be flowing forwards.

In addition, much like $\langle\Delta\mathcal{S}^{y}_{\mathrm{tot}}\rangle$, $\langle\Delta\mathcal{S}^{y|x}_{\mathrm{info}}\rangle$ has no bounds on its sign: the information transfer in the reverse time direction *can* be greater than in the forward direction. However, we argue that the total, or composite, arrow of time, formed from both the physical and information theoretic components ought to hold in all circumstances. This total arrow of time, in a heuristic

sense, modifies its usual conception so as to contrast not only the likelihood of a path with its time reversed counterpart, but the likelihood of a path in conjunction with how much information can be or is extracted from it. Instead of contrasting the functional $\hat{p}[\mathbf{y}^{\tau}_{t_0}|\{\mathbf{x}^{\tau}_{t_0}\}]$ in the forward and reverse time directions with one expected to be larger in the forward time direction, we instead use $\hat{p}[\mathbf{y}^{\tau}_{t_0}|\{\mathbf{x}^{\tau}_{t_0}\}]\exp[I_m(\mathbf{x}_{t_0}, \mathbf{y}_{t_0}) + \mathcal{T}_{y\to x}]$. The demonstration of the positivity of the total arrow of time is straightforward from the bipartite structure utilized which leads to

$$\Delta\mathcal{S}^{y}_{\mathrm{tot}} + \Delta\mathcal{S}^{y|x}_{\mathrm{info}} = \ln\frac{p[\mathbf{x}^{\tau}_{t_0}, \mathbf{y}^{\tau}_{t_0}]}{p[\mathbf{x}^{t_0}_{\tau}, \mathbf{y}^{t_0}_{\tau}]} - \ln\frac{p[\mathbf{x}^{\tau}_{t_0}]}{p[\mathbf{x}^{t_0}_{\tau}]}$$
$$= \Delta\mathcal{S}_{\mathrm{tot}} - \Delta\mathbb{S}^{x}_{\mathrm{tot}}. \tag{39}$$

Such a quantity is known to obey an integral fluctuation theorem [63], demonstrated in Appendix C, and is therefore positive in the mean. However, whilst the non-Markov structure can make the final quantity in eq. (39) ambiguous when interpreted as a physical entropy production, the transfer of information is well defined, explicitly with non-Markovian dependence, giving such a quantity a concrete interpretation in terms of an arrow of time comprising both physical and information flows forward and backwards in time.

At this point we recognize that both the numerator and denominator in the information theoretic arrow of time are precisely the additional encoding cost identified in eq. (30) such that the composite arrow of time can be written

$$\Delta\mathcal{S}^{y}_{\mathrm{tot}} + \Delta\mathcal{S}^{y|x}_{\mathrm{info}} = \ln\frac{p(\mathbf{y}_{t_0})\hat{p}[\mathbf{y}^{\tau}_{t_0}|\mathbf{y}_{t_0}, \{\mathbf{x}\}^{\tau}_{t_0}]}{p(\mathbf{y}_{\tau})\hat{p}[\mathbf{y}^{t_0}_{\tau}|\mathbf{y}_{\tau}, \{\mathbf{x}\}^{t_0}_{\tau}]}$$
$$+ d(p_{\mathbf{y}|\mathbf{x}}||\hat{p}_{\mathbf{y}|\{\mathbf{x}\}}) - d(p_{\mathbf{y}^{\dagger}|\mathbf{x}^{\dagger}}||\hat{p}_{\mathbf{y}^{\dagger}|\{\mathbf{x}^{\dagger}\}}). \tag{40}$$

The addition of this additional encoding cost can be seen as a re-weighting of the measure that considers $\mathbf{x}$ to be a protocol, to the measure which aligns with the actual observed conditional probability such that $p_w[\mathbf{y}^{\tau}_{t_0}|\{\mathbf{x}\}^{\tau}_{t_0}] = \hat{p}[\mathbf{y}^{\tau}_{t_0}|\{\mathbf{x}\}^{\tau}_{t_0}]\exp[I_m(\mathbf{x}_{t_0}, \mathbf{y}_{t_0}) + \mathcal{T}_{y\to x}] = p[\mathbf{y}^{\tau}_{t_0}|\mathbf{x}^{\tau}_{t_0}]$. As such the composite arrow of time can be written as an integral over Kullback-Leibler divergences in contrast to eq. (37)

$$\left\langle\Delta\mathcal{S}^{y}_{\mathrm{tot}} + \Delta\mathcal{S}^{y|x}_{\mathrm{info}}\right\rangle$$
$$= \int d\mathbf{x}^{\tau}_{t_0}\int d\mathbf{y}^{\tau}_{t_0}p[\mathbf{y}^{\tau}_{t_0}, \mathbf{y}^{\tau}_{t_0}]\ln\frac{p_w[\mathbf{y}^{\tau}_{t_0}|\{\mathbf{x}\}^{\tau}_{t_0}]}{p_w[\mathbf{y}^{t_0}_{\tau}|\{\mathbf{x}\}^{t_0}_{\tau}]}$$
$$= \int d\mathbf{x}^{\tau}_{t_0}p[\mathbf{x}^{\tau}_{t_0}]\int d\mathbf{y}^{\tau}_{t_0}p[\mathbf{y}^{\tau}_{t_0}|\mathbf{x}^{\tau}_{t_0}]\ln\frac{p[\mathbf{y}^{\tau}_{t_0}|\mathbf{x}^{\tau}_{t_0}]}{p[\mathbf{y}^{t_0}_{\tau}|\mathbf{x}^{\tau}_{t_0}]}$$
$$= \left\langle D_{KL}(p_{\mathbf{y}|\mathbf{x}}||p^{\dagger}_{\mathbf{y}|\mathbf{x}})\right\rangle_{\mathbf{x}} \geq 0. \tag{41}$$

In contrast to the previous path integrals, the notation $\langle\rangle_x$ indicates a path integral over all trajectories in $\mathbf{x}$

only. Furthermore, by considering the quantities related to the forward time direction we interpret the various quantities as contributions to the encoding cost of the specific realization $\mathbf{y}_{t_0}^\tau$ given $\mathbf{x}_{t_0}^\tau$. Under this interpretation $-\ln[\hat{p}[\mathbf{y}_{t_0}^\tau|\{\mathbf{x}\}_{t_0}^\tau]]$ is the code length of $\mathbf{y}_{t_0}^\tau$ under the non-optimal encoding strategy which considers $\mathbf{x}_{t_0}^\tau$ as a protocol, $d(p_{\mathbf{y}|\mathbf{x}}||\hat{p}_{\mathbf{y}|\{\mathbf{x}\}})$ as before is the additional cost associated with utilizing a non-optimal encoding strategy and $-\ln[p[\mathbf{y}_{t_0}^\tau|\mathbf{x}_{t_0}^\tau]]$ is the code length under the optimal encoding. When averaged over all realizations of $\mathbf{x}_{t_0}^\tau$ we can write the mean encoding costs under the non-optimal measure as the average cross entropy, $H_c(p(y), q(y)) = -\int dy p(y) \ln q(y)$,

$$\langle H_c(p_{\mathbf{y}|\mathbf{x}}, \hat{p}_{\mathbf{y}|\{\mathbf{x}\}})\rangle_\mathbf{x} = \langle H(p_{\mathbf{y}|\mathbf{x}})\rangle_\mathbf{x} + \langle D_{KL}(p_{\mathbf{y}|\mathbf{x}}||\hat{p}_{\mathbf{y}|\{\mathbf{x}\}})\rangle_\mathbf{x} \tag{42}$$

and where $H(p(y)) = -\int dy p(y) \ln p(y)$. In short, the information transfer which forms $\Delta\mathcal{S}_{\text{info}}^{y|x}$ can be viewed as a statistical re-weighting which changes the naive measure, now on $\mathbf{y}$, into the true measure. With this inclusion the total arrow of time is then formed from code lengths under the true measure and leads to the positivity of the composite arrow of time as there exists no other measure which produces, on average, a superior encoding scheme.

Further, the form in eq. (42) suggests an interesting interpretation in terms of the principle of minimum cross entropy or principle of minimum discrimination [72–74]. It concerns situations where there is a given reference measure $\hat{p}$ which is approximating some unknown measure $p$ and an improved measure $p'$ is to be chosen given some new evidence, usually in the form of some constraint(s). The principle then states that $p'$ should be chosen such that it minimizes the discrimination entropy (the Kullback-Leibler divergence) or equivalently, the cross entropy. Such a principle can be seen as an extension of the principle of maximum entropy. The relevance here, we argue, is that the initial reference measure, given knowledge of the dynamics of $\mathbf{y}$, but not of $\mathbf{x}$, is sensibly chosen to be the naive measure. But as continued measurements are made evidence might arrive that this measure is not the true measure. We are then faced with the inference problem of determining the true measure given this new evidence. The principle of minimum discrimination then suggests that the best approximation to the true measure is that which minimizes the Kullback-Leibler divergence between the naive measure and the new approximation. However, because of eq. (29) we can interpret this as the measure which leads to the least information transfer from $\mathbf{y}$ to $\mathbf{x}$: the best approximation of the true measure is through an estimation of the dynamics of $\mathbf{x}$ which can reproduce the necessary constraints with the *minimum* information about $\mathbf{y}$ as possible.

Finally, since the form of $\Delta\mathcal{S}_{\text{info}}^{y|x}$ is involutive ($-\Delta\mathcal{S}_{\text{info}}^{y|x}$ is the value of $\Delta\mathcal{S}_{\text{info}}^{y|x}$ when time is reversed) we can

measure this one quantity in both time directions to infer the arrow of time in the same way one usually measures $\Delta\mathcal{S}_{\text{tot}}^y$. Alternatively, we can consider attempting to encode the forward and time reversed behaviour with the implication that the shortest code length is associated with the forward time direction. Encoding under the naive measure reflects the entropy production, but may mislead on average. Encoding under the optimal scheme will not mislead on average and the corrective term is the exported information given by the initial mutual information and the pathwise transfer entropy.

Lastly, we point out that if we restrict ourselves to the steady state, such that $\lambda = \text{const}$, and assume a finite correlation of the stochastic behaviour of the coarse grained system with its past, by considering a process of infinite duration such that the contributions near the time origin and horizon become negligible, it follows that

$$\Delta\dot{S}_i^y + T_{y\to x} - T_{y\to x}^\dagger \geq 0 \tag{43}$$

where

$$T_{y\to x} = \lim_{t\to\infty} \frac{1}{t}\langle\mathcal{T}_{y\to x}\rangle$$
$$T_{y\to x}^\dagger = \lim_{t\to\infty} \frac{1}{t}\langle\mathcal{T}_{y\to x}^\dagger\rangle. \tag{44}$$

Such a relation is expected to provide a stronger bound on the entropy production $\Delta\dot{S}_i^y$ near equilibrium in the total system than just the transfer entropy rate alone [31] since the information transfer in the reverse process, is a non vanishing quantity equal in magnitude to the forward information transfer at equilibrium. We also see that if only subsystem $y$ is in thermodynamic equilibrium we have $T_{y\to x} \geq T_{y\to x}^\dagger$ such that information can always be extracted more effectively running forwards in time versus backwards in time, but if the measuring apparatus is also at equilibrium such that the total device is in equilibrium no distinction can be made and we have $T_{y\to x} = T_{y\to x}^\dagger$.

## VI. CONNECTION WITH PREVIOUS RESULTS

The result we have presented is, as written, specifically formulated for systems undergoing so called autonomous feedback; systems where the natural coupling and resulting dynamics can be seen to lead to information exchange between identifiable components of the total device. It is, however, instructive to examine how such a result can be seen to lead to previous results by imagining it for different setups by making further assumptions about the nature of time reversal.

In contrast to the setup considered here is the setup known as non-autonomous feedback. In such setups rather than describing two components of a joint system one describes a single system of which measurements

are made by an undescribed controller which then uses the information gained by the measurement to adjust the protocol. In these setups system $\mathbf{x}$ is measured by a measurement device recording outcomes in a memory controller $\mathbf{y}$ which then changes the protocol $\lambda(\mathbf{y})$ [27]. Whilst it is instructive for the consideration of the thermodynamics of the system to treat it in this way, the dynamics can still be entirely captured through a combined *bipartite* system whereby $\mathbf{y}$ plays the dual role of measurement device and protocol. This is because, in line with the fundamentals of stochastic thermodynamics, protocols are stationary whilst the system transitions and any measurement results are statistically independent of the future states of the system.

Once such a connection is made the distinction between such systems is in the nature of time reversal: in autonomous feedback systems both subsystems exhibit dependence on one another regardless of time reversal, whilst for non-autonomous feedback systems the joint measurement/protocol variable becomes independent of the system variable upon time reversal [75]. Such an effect may be imagined by considering the usual bipartite system with a constant protocol variable (for the whole system) $\lambda$, such as a constant magnetic field, which transforms upon time reversal to a value $\lambda^\dagger$ in such a way that mediates the ability of $\mathbf{y}$ to sense $\mathbf{x}$. In this set up $\mathbf{x}$ still depends on $\mathbf{y}$ in both the forward and reverse processes leaving the entropy production of $\mathbf{x}$ unchanged as well as the initial and final mutual informations and forward transfer entropy. However, the reverse transfer entropy now explicitly vanishes for all realizations. Consequently the information theoretic arrow of time reduces to

$$\Delta \mathcal{S}_{\text{info}}^{y|x} = I_m(\mathbf{x}_{t_0}, \mathbf{y}_{t_0}) + \mathcal{T}_{\mathbf{y} \to \mathbf{x}} - I_m(\mathbf{x}_\tau, \mathbf{y}_\tau). \quad (45)$$

Consequently eq. (39), which obeys an integral fluctuation relation, reduces to

$$\Delta \mathcal{S}_{\text{tot}}^y + \Delta \mathcal{S}_{\text{info}}^{y|x} = \Delta \mathcal{S}_i^y + I_m(\mathbf{x}_{t_0}, \mathbf{y}_{t_0}) + \mathcal{T}_{\mathbf{y} \to \mathbf{x}} - I_m(\mathbf{x}_\tau, \mathbf{y}_\tau) \quad (46)$$

thus precisely recovering the result in eqs. (5) & (6) in [71]. Written as a Kullback-Leibler divergence analogous to eq. (41) the average of such a quantity is given by

$$\langle \Delta \mathcal{S}_{\text{tot}}^y + I_m(\mathbf{x}_{t_0}, \mathbf{y}_{t_0}) + \mathcal{T}_{\mathbf{y} \to \mathbf{x}} - I_m(\mathbf{x}_\tau, \mathbf{y}_\tau) \rangle$$
$$= \int d\mathbf{x}_{t_0}^\tau p(\mathbf{x}_{t_0}^\tau)$$
$$\times \int d\mathbf{y}_{t_0}^\tau p[\mathbf{y}_{t_0}^\tau | \mathbf{x}_{t_0}^\tau] \ln \frac{p[\mathbf{y}_{t_0}^\tau | \mathbf{x}_{t_0}^\tau]}{p(\mathbf{y}_\tau | \mathbf{x}_\tau) \hat{p}[\mathbf{y}_\tau^{t_0} | \mathbf{y}_\tau, \{\mathbf{x}\}_\tau^{t_0}]}$$
$$\geq 0. \quad (47)$$

Further, if we insist that any such process, forward or backwards in time, is first prepared so that the system and the controller/protocol variable are uncorrelated the mutual information terms vanish and we find

$$\Delta \mathcal{S}_{\text{tot}}^y + \Delta \mathcal{S}_{\text{info}}^{y|x} = \Delta \mathcal{S}_{\text{tot}}^y + \mathcal{T}_{\mathbf{y} \to \mathbf{x}} \quad (48)$$

recovering the well known result found in eq. (6) in [24] and eq. (65) in [27]. Similarly, in the stationary state the mean mutual information contributions cancel and we confirm the result in eq. (37) in [29].

Finally we can make a connection with the result in [76–78] by noticing that if the normalization terms within such descriptions are identified as $Z_1 = (p(x_i|x_{i+1}))^{-1}$ and $Z_2 = (p(x_i|x_{i+1}, y_{i+1}))^{-1}$, the corrective log ratio term, $\ln(Z_1/Z_2)$, becomes the reverse transfer entropy. Consequently the contained $\Delta S_{ext}(x)$ should be recognized as the difference between the exported entropy production attributed to $\mathbf{x}$ at different levels of description.

## VII.   EXAMPLE SYSTEM

Relevant systems for such results are numerous [30], however beyond numerical computation of the constituent quantities, exact analytical results are challenging to find. An instructive example where an analytical solution can be found in an appropriate limit is that of two simple harmonic oscillators, coupled by a harmonic potential, subject to distinct environments at different temperatures described by over-damped Langevin equations. In this system each oscillator represents one of the two sub-systems $x = \mathbf{x}$ and $y = \mathbf{y}$ in the preceding development with the coupling allowing us to frame the dynamics in terms of autonomous feedback when we take particular observer perspectives. By altering the temperatures that the two oscillators are exposed to we equivalently change the statistical properties of the feedback and cool/heat the oscillators relative to their respective heat baths. The equations of motion for the system take the form of the following stochastic differential equations

$$dx = -\frac{1}{\gamma_x} kx \, dt - \frac{1}{\gamma_x} k_c(x - y) \, dt + \sqrt{\frac{2k_B T_x}{\gamma_x}} dW^x$$

$$dy = -\frac{1}{\gamma_y} ky \, dt - \frac{1}{\gamma_y} k_c(y - x) \, dt + \sqrt{\frac{2k_B T_y}{\gamma_y}} dW^y \quad (49)$$

where $\gamma_x, \gamma_y$ are the damping coefficients for $x, y$ respectively, $k$ the spring constant of both oscillators, $k_c$ the spring constant of the coupling potential, $T_x, T_y$ the temperatures of the environments the oscillators are exposed to and $dW_x$ & $dW_y$ are two uncorrelated Wiener processes simulating the thermal environments. Such a system is continuous in time and bipartite since the current can readily be divided into two components because of this lack of correlation between the Wiener processes [31]. Whilst such systems permit solutions to the transfer entropy rate for stationary processes in the limit $s \to \infty$ [31, 61], we are not aware of such a solution for the reverse transfer entropy rate. Consequently, for consistency, we approximate the transfer and reverse transfer entropies using a method involving small time propagators, effectively calculating the transfer entropy rates

with $\lim s \to 0$, (See appendix D), quantites which have been investigated in [61]. These then become accurate, such that $T_{y \to x} = \dot{\mathcal{T}}_{y \to x}^{(0)}$, when the source variable, $y$, is much faster than the target, $x$. This occurs when $\gamma_y << \gamma_x$. In such a limit we find

$$\frac{d\langle \Delta \mathcal{S}_{\text{tot}}^y \rangle}{dt} = \frac{k_c^2 (T_x - T_y)}{(k + k_c)\gamma_x T_y} + O\left(\frac{\gamma_y}{\gamma_x}\right)$$

$$\frac{d\langle \mathcal{T}_{y \to x}^{(0)} \rangle}{dt} = \frac{k_c^2 T_y}{4(k + k_c)\gamma_x T_x} + O\left(\frac{\gamma_y}{\gamma_x}\right)$$

$$= \lim_{\gamma_y/\gamma_x \to 0} T_{y \to x}$$

$$\frac{d\langle \mathcal{T}_{y \to x}^{\dagger,(0)} \rangle}{dt} = \frac{k_c^2 (4T_x - 3T_y)}{4(k + k_c)\gamma_x T_x} + O\left(\frac{\gamma_y}{\gamma_x}\right)$$

$$= \lim_{\gamma_y/\gamma_x \to 0} T_{y \to x}^{\dagger}. \qquad (50)$$

Such results are illustrated in fig. (1). When $T_y > T_x$ the coupling, when viewed from the perspective of $x$ as a protocol, serves to cool the particle and results in an apparent negative entropy production in $y$. This is offset by a positive transfer entropy from $x$ to $y$ providing a familiar bound on their sum [27, 71]. However, the transfer entropy is postive for *all* $T_y$. Alternatively, if we examine the information theoretic arrow of time we see that this is positive *only* in the regime $T_y > T_x$. Consequently when $T_y = T_x$, when the joint system is in equilibrium, we recover no discernable arrow of time whether that be thermodynamic or information theoretic, $\langle \dot{\mathcal{S}}_{\text{tot}}^y \rangle = \langle \dot{\mathcal{S}}_{\text{tot}}^x \rangle = \langle \dot{\mathcal{S}}_{\text{info}}^{y|x} \rangle = \langle \dot{\mathcal{S}}_{\text{info}}^{x|y} \rangle = 0$. Interestingly, in the regime $0 \le T_y \le 4T_x/3$, the reverse measure is not too far from the forward measure such that the reverse transfer entropy is positive: the source $y$ still provides information about the destination $x$ for the forward time direction whilst time is played in reverse. Within this bound the information theoretic arrow of time provides a stronger bound than the transfer entropy alone and for $T_x < T_y < 4T_x/3$ this occurs whilst $\langle \Delta \dot{\mathcal{S}}_{\text{tot}}^y \rangle$ is negative. Further, when $T_x > T_y$ such that $\langle \Delta \dot{\mathcal{S}}_{\text{tot}}^y \rangle$ is positive the reverse transfer entropy is positive and larger than the transfer entropy suggesting that when the reversed time measure over samples less dissipative paths in the source, more information can be transferred from that source.

## VIII. DISCUSSION AND CONCLUSION

In this paper we have introduced a new quantity we call the reverse transfer entropy in the hope of creating deeper connections between information theory and thermodynamics. We have then described such a quantity for bipartite systems and related it to the usual transfer entropy and the observed irreversibility of a system at two different levels of description which differ in their knowledge of the source variable. Further, by making a direct comparison between physical flows
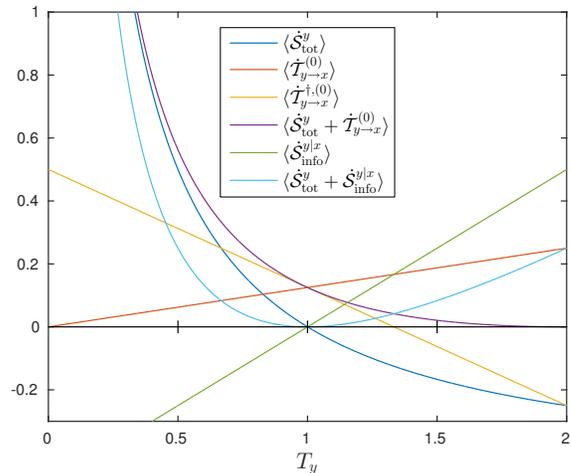


FIG. 1. Entropy production, transfer & reverse transfer entropy, and information theroretic arrow of time rates for the coupled simple harmonic oscillators at two temperatures for values $k_B = T_x = \gamma_x = k = k_c = 1$ in the limit $\lim_{\gamma_y/\gamma_x \to 0}$.

represented by observed probabilities and information flows represented by the transfer entropy we have constructed an information theoretic arrow of time. This information theoretic arrow of time can, on its own, either inform or mislead about the true direction of the arrow of time. But when paired with the physical arrow of time, which similarly can mislead in the presence of feedback, we have presented a quantity, comprising both the probability of a realization and the information which it can export, which reliably informs about the arrow of time with or without feedback.

By considering different models for time reversal we have then made connections with previously known results in the literature. We argue that by doing so we have illustrated that it can be viewed that it is this composite arrow of time which underlies many of the results involving feedback as they are recovered by introducing greater restrictions and assumptions about time reversal which act to set certain quantities in the composite arrow of time, namely the reverse transfer entropy, to zero. This difference in perspective is subtle, but significant: rather than viewing bounds on the entropy production as emanating from a physical probabilities of the forward and reverse paths alone such that information theoretic quantities can be derived, we start by placing an equal weighting on the information theoretic quantities themselves and view them as relevant to the combined thermodynamics. We have also presented a succinct alternative interpretation of the arrow of time, valid in any system: it is the ability, using a single encoding strategy, to reliably encode the forward time behaviour more efficiently than the reverse. But this only holds when the optimal encoding scheme is utilized. If one uses a scheme based on the naive

measure, it leads to the entropy production (eqs. (36) & (37)), but may reach the wrong conclusion about the flow of time. The additional cost of this error is precisely the information exported by the system in addition to a boundary given by the initial mutual information. Such an interpretation further highlights the difference between the approach outlined here and previous results [29, 31, 71], because the composite arrow of time uses the same encoding scheme in both time directions whereas previous results that do not include the reverse transfer entropy, when interpreted as a comparison between encoding costs, entail a change in scheme in the reverse time direction (eq. (47)). We anticipate further richness in the comparison with other information theoretic quantities with those in stochastic thermodynamics.

## Appendix A: Incorporation of mechanisms into transfer entropy

To fully describe transfer entropy in the context of stochastic physical models we must acknowledge that a given transition may occur by one of several mechanisms $\nu$ [40–46], the existence of which additively contribute to the probability of a transition ($p(x'|x) = \sum_\nu p_\nu(x'|x)$). An example might be a system connected to multiple heat or particle reservoirs which each provide a separable source of noise and can thus both be considered in isolation as well as together. This extension allows the development of the notion of a transition being identifiably induced by one of such sources thus providing additional information about any transition. The incorporation of this additional knowledge of *how* a transition occurs can be incorporated into the transfer of information through an expansion of the state space into a pseudo state-mechanism space, $(\mathbf{x}_i, \mathbf{y}_i) \to (\mathbf{x}_i, \mathbf{y}_i, \nu_i)$, whereby $\nu_n$ indicates the mechanism utilized in the transition from state $\mathbf{x}_{n-1} \to \mathbf{x}_n$.

The key feature of introducing additional mechanisms into physical systems is that the rates from the contributing mechanisms are additive such that

$$W(\mathbf{x}, \mathbf{y}|\mathbf{x}', \mathbf{y}') = \sum_\nu W_\nu(\mathbf{x}, \mathbf{y}|\mathbf{x}', \mathbf{y}'). \qquad (A1)$$

Considering, for simplicity, continuous Markov processes described by a master equation

$$\frac{\partial P(\mathbf{x}, \mathbf{y})}{\partial t} = \sum_{\mathbf{x}', \mathbf{y}'} W(\mathbf{x}, \mathbf{y}|\mathbf{x}', \mathbf{y}') P(\mathbf{x}', \mathbf{y}'), \qquad (A2)$$

we can incorporate the additional knowledge of the mechanisms used into the transfer entropy by building a pseudo state space which adds the mechanism of the transition into that state so that $P(\mathbf{x}, \mathbf{y}) \to P(\mathbf{x}, \mathbf{y}, \nu)$. Consequently the transition rates are now

of the form $W(\mathbf{x}, \mathbf{y}, \nu|\mathbf{x}', \mathbf{y}', \nu')$. However, the rates cannot depend on the mechanism of the previous transition so we write $W(\mathbf{x}, \mathbf{y}, \nu|\mathbf{x}', \mathbf{y}', \nu') = W(\mathbf{x}, \mathbf{y}, \nu|\mathbf{x}', \mathbf{y}') = W_\nu(\mathbf{x}, \mathbf{y}|\mathbf{x}', \mathbf{y}')$. The master equation then becomes

$$\frac{\partial P(\mathbf{x}, \mathbf{y}, \nu)}{\partial t} = \sum_{\mathbf{x}', \mathbf{y}', \nu'} W(\mathbf{x}, \mathbf{y}, \nu|\mathbf{x}', \mathbf{y}') P(\mathbf{x}', \mathbf{y}', \nu')$$
$$= \sum_{\mathbf{x}', \mathbf{y}'} W_\nu(\mathbf{x}, \mathbf{y}|\mathbf{x}', \mathbf{y}') P(\mathbf{x}', \mathbf{y}'). \qquad (A3)$$

Summing over $\nu$ gives the original master equation ensuring such a description is consistent with the original system.

The transfer entropy is then formally the transfer $T_{y \to x, \nu}^{(s,r)}$ on the pseudo space such that, taking for simplicity $s = r = 0$,

$$T_{y \to x, \nu}^{(0,0)} = \sum_{\mathbf{x}, \mathbf{x}', \mathbf{y}, \mathbf{y}', \nu, \nu'} P(\mathbf{x}', \mathbf{y}', \nu') W(\mathbf{x}, \mathbf{y}, \nu|\mathbf{x}', \mathbf{y}')$$
$$\times \ln \frac{W(\mathbf{x}, \nu|\mathbf{x}', \mathbf{y}')}{W(\mathbf{x}, \nu|\mathbf{x}')}$$
$$= \sum_{\mathbf{x}, \mathbf{x}', \mathbf{y}', \nu} P(\mathbf{x}', \mathbf{y}') W_\nu(\mathbf{x}|\mathbf{x}', \mathbf{y}') \ln \frac{W_\nu(\mathbf{x}|\mathbf{x}', \mathbf{y}')}{W_\nu(\mathbf{x}|\mathbf{x}')}$$
$$= \sum_\nu T_{y \to x}^{(0,0), \nu} \qquad (A4)$$

where $T_{y \to x}^{(0,0), \nu}$ is the transfer entropy of the original system if all transitions were assumed to be caused by that configuration of mechanisms $\nu$. From the log sum inequality it follows that

$$T_{y \to x, \nu}^{(0,0)} = \sum_{\mathbf{x}, \mathbf{x}', \mathbf{y}', \nu} P(\mathbf{x}', \mathbf{y}') W_\nu(\mathbf{x}|\mathbf{x}', \mathbf{y}') \ln \frac{W_\nu(\mathbf{x}|\mathbf{x}', \mathbf{y}')}{W_\nu(\mathbf{x}|\mathbf{x}')}$$
$$\geq \sum_{\mathbf{x}, \mathbf{x}', \mathbf{y}'} P(\mathbf{x}', \mathbf{y}') \sum_\nu W_\nu(\mathbf{x}|\mathbf{x}', \mathbf{y}') \ln \frac{\sum_\nu W_\nu(\mathbf{x}|\mathbf{x}', \mathbf{y}')}{\sum_\nu W_\nu(\mathbf{x}|\mathbf{x}')}$$
$$= \sum_{\mathbf{x}, \mathbf{x}', \mathbf{y}'} P(\mathbf{x}', \mathbf{y}') W(\mathbf{x}|\mathbf{x}', \mathbf{y}') \ln \frac{W(\mathbf{x}|\mathbf{x}', \mathbf{y}')}{W(\mathbf{x}|\mathbf{x}')}$$
$$= \bar{T}_{y \to x}^{(0,0)} \qquad (A5)$$

where the last quantity is the transfer entropy calculated without knowledge of the mechanisms in the system. Incorporating these features along with the inclusion of the protocol transforms the transfer entropy in physical systems into a transfer entropy, conditional upon the protocol, $\lambda$, from variable $\mathbf{y}$ to the joint states on the state mechanism pseudo space $(\mathbf{x}, \nu)$ and is of the form

$$T_{\mathbf{y} \to \mathbf{x}, \nu|\lambda_{n-m}^{n-1}}^{(k,l,m)}(n) = \left\langle \ln \frac{p(\mathbf{x}_n, \nu_n|\mathbf{x}_{n-k}^{n-1}, \mathbf{y}_{n-l}^{n-1}, \lambda_{n-m}^{n-1})}{p(\mathbf{x}_n, \nu_n|\mathbf{x}_{n-k}^{n-1}, \lambda_{n-m}^{n-1})} \right\rangle. \qquad (A6)$$

Explicitly, this is

$$
\begin{aligned}
T^{(k,l,m)}_{\mathbf{y}\to\mathbf{x},\nu|\lambda^{n-1}_{n-m}}(n) &= \sum_{\nu_n}\sum_{\mathbf{x}^n_{n-k}}\sum_{\mathbf{y}^{n-1}_{n-l}} p(\mathbf{x}^n_{n-k},\mathbf{y}^{n-1}_{n-l},\nu_n|\lambda^{n-1}_{n-m}) \\
&\quad \times \ln \frac{p(\mathbf{x}_n,\nu_n|\mathbf{x}^{n-1}_{n-k},\mathbf{y}^{n-1}_{n-l},\lambda^{n-1}_{n-m})}{p(\mathbf{x}_n,\nu_n|\mathbf{x}^{n-1}_{n-k},\lambda^{n-1}_{n-m})} \\
&= \sum_{\nu_n} T^{(k,l,m),\nu_n}_{\mathbf{y}\to\mathbf{x}|\lambda^{n-1}_{n-m}}(n) \qquad (A7)
\end{aligned}
$$

where the final sum is over transfer entropies for the same system/observations with all transitions assumed to occur according to mechanism $\nu_n$. We note that by the log sum inequality we have

$$
T^{(k,l,m)}_{\mathbf{y}\to\mathbf{x},\nu|\lambda^{n-1}_{n-m}}(n) \geq \bar{T}^{(k,l,m)}_{\mathbf{y}\to\mathbf{x}|\lambda^{n-1}_{n-m}}(n). \qquad (A8)
$$

Here, $\bar{T}$, is the information transfer in the absence of knowledge of $\nu$. This indicates, much like the entropy production in [45], that a failure to account for the additional knowledge of the mechanisms will underestimate the transfer of information.

## Appendix B: Transfer entropy and reverse transfer entropy for continuous time systems

Here we start from eq. (6) and recognize that either the numerator or denominator in such a quantity, for a path running from $t = t_0$ to $t = t_{N+1}$ consisting of $N_x$ transitions in $\mathbf{x}$ with the $i$th transition occurring at time $t_i$ being into state $\mathbf{x}_i$, is represented by

$$
\begin{aligned}
&P(\mathbf{x}=\mathbf{x}_0 \forall t \in [t_0,t_1]) \times \\
&\prod_{i=1}^{N_x} P(\mathbf{x}_{i-1}\to\mathbf{x}_i, t \in dt_i)P(\mathbf{x}=\mathbf{x}_i\forall t \in [t_i,t_{i+1}]).
\end{aligned}
$$
$$(B1)$$

Transition terms are determined by instantaneous transition rates such that

$$
P(\mathbf{x}_{i-1}\to\mathbf{x}_i, t \in dt_i) = W_{\mathbf{x}_{i-1}\mathbf{x}_i}dt_i \qquad (B2)
$$

where $W_{\mathbf{x}_{i-1}\mathbf{x}_i} = W[\mathbf{x}_i|\mathbf{x}^{t_i}_{t_i-s}]$, such that $\mathbf{x}(t_i) = \lim_{t\to t_i-}\mathbf{x}(t) = \mathbf{x}_{i-1}$, for the coarse grained path probability appearing in the denominator and $W_{\mathbf{x}_{i-1}\mathbf{x}_i} = W[\mathbf{x}_i|\mathbf{x}^{t_i}_{t_i-s},\mathbf{y}^{t_i}_{t_i-s}] = \sum_{\mathbf{y}'} W[\mathbf{x}_i,\mathbf{y}'|\mathbf{x}^{t_i}_{t_i-s},\mathbf{y}^{t_i}_{t_i-r}]$ for the conditional process appearing in the numerator. In each case the stationary probabilities can be described by the

limit of the sequence

$$
\begin{aligned}
P(\mathbf{x}=\mathbf{x}_i\forall t \in [t_i,t_{i+1}]) &= \lim_{dt\to 0}\prod_{t_j=t_i}^{t_j=t_{i+1}}\left(1-\sum_{\mathbf{x}'\neq\mathbf{x}_j}W_{\mathbf{x}_j\mathbf{x}'}dt\right) \\
&= \lim_{dt\to 0}\exp\left[-\sum_{t_j=t_i}^{t_j=t_{i+1}}\sum_{\mathbf{x}'\neq\mathbf{x}_j}W_{\mathbf{x}_j\mathbf{x}'}dt\right] \\
&= \exp\left[-\int_{t_i}^{t_{i+1}}dt\left(\sum_{\mathbf{x}'\neq\mathbf{x}(t)}W_{\mathbf{x}(t)\mathbf{x}'}\right)\right] \\
&= \exp\left[-\int_{t_i}^{t_{i+1}}dt\,\kappa_{\mathbf{x}(t)}\right] \qquad (B3)
\end{aligned}
$$

where it should be noted that, despite the exponential form, the Markov property is not retained due to the integral over the mean escape rate which can explicitly change while the system is not transitioning even for time homogeneous processes since the history encoded into the transition rates changes as time progresses. Simply rearranging eq. (B1) and taking the logarithm gives

$$
\ln\left(\prod_{i=1}^{N_x}W_{\mathbf{x}_{i-1}\mathbf{x}_i}dt\right) - \int_{t_0}^{t_{N+1}}\kappa_{\mathbf{x}(t)}dt \qquad (B4)
$$

being the logarithm of a path probability density with dimension $dt^{N_x}$. Inserting the appropriate transition rates and taking the difference between the conditional and marginal processes then gives eq. (6).

To construct the reverse transfer entropy we consider the time reversed path, assuming even variables, given by $\mathbf{x}^\dagger(t) = \mathbf{x}(t_{N_x+1}+t_0-t)$ so that $\mathbf{x}^\dagger_i = \mathbf{x}_{N_x-i}$. The reverse form of eq. (B4) is then given by

$$
\ln\left(\prod_{i=1}^{N_x}W^*_{\mathbf{x}^\dagger_{i-1}\mathbf{x}^\dagger_i}dt\right) - \int_{t_0}^{t_{N+1}}\kappa^*_{\mathbf{x}^\dagger(t)}dt \qquad (B5)
$$

where for the conditional process

$$
\kappa^*_{\mathbf{x}^\dagger(t)} = \sum_{\mathbf{x}'\neq\mathbf{x}^\dagger(t)}\sum_{\mathbf{y}'}W(\mathbf{x}',\mathbf{y}'|\mathbf{x}^{\dagger t}_{t-s},\mathbf{y}^{\dagger t}_{t-r}) \qquad (B6)
$$

with the coarse grained process being analogous. Substituting in with the definition of the reversed path and recognizing that the path history runs in the opposite direction we have

$$
\begin{aligned}
\kappa^*_{\mathbf{x}^\dagger(t)} &= \kappa^*_{\mathbf{x}(t_{N_x+1}+t_0-t)} \\
&= \sum_{\mathbf{x}'\neq\mathbf{x}(t_{N_x+1}+t_0-t)}W^*_{\mathbf{x}_{t_{N_x+1}+t_0-t}\mathbf{x}'} \\
&= \sum_{\mathbf{x}'\neq\mathbf{x}(t_{N_x+1}+t_0-t)}\sum_{\mathbf{y}'} \\
&\quad W(\mathbf{x}',\mathbf{y}'|\mathbf{x}^{t_{N_x+1}+t_0-t}_{t_{N_x+1}+t_0-t+s},\mathbf{y}^{t_{N_x+1}+t_0-t}_{t_{N_x+1}+t_0-t+r}) \qquad (B7)
\end{aligned}
$$

with the direction of the histories distinguishing $\kappa^*$ and $W^*$ from $\kappa$ and $W$. We thus arrive at

$$\ln\left(\prod_{i=1}^{N_x} W^*_{\mathbf{x}_{N_x-i+1}\mathbf{x}_{N_x-i}}dt\right) - \int_{t_0}^{t_{N+1}} \kappa^*_{\mathbf{x}(t_{N_x+1}+t_0-t)}dt. \tag{B8}$$

Changing to the time variable $t' = t_{N_x+1} + t_0 - t$ which runs in the same direction as the forward process such that $t_i \to t_{N_x-i+1}$ gives

$$\ln\left(\prod_{i=1}^{N_x} W^*_{\mathbf{x}_{N_x-i+1}\mathbf{x}_{N_x-i}}dt\right) - \int_{t_0}^{t_{N+1}} \kappa^*_{\mathbf{x}(t')}dt'$$

$$= \ln\left(\prod_{i=1}^{N_x} W^*_{\mathbf{x}_i\mathbf{x}_{i-1}}dt\right) - \int_{t_0}^{t_{N+1}} \kappa^*_{\mathbf{x}(t)}dt \tag{B9}$$

which leads to the form of the reverse transfer entropy in eq. (16). We note that when the process is Markov we have $W^* = W$ and $\kappa^* = \kappa$.

## Appendix C: Demonstration of the integral fluctuation theorem for the composite arrow of time

We have

$$\Delta\mathcal{S}^y_{\text{tot}} + \Delta\mathcal{S}^{y|x}_{\text{info}}$$
$$= \ln\frac{p(\mathbf{y}_{t_0})\hat{p}[\mathbf{y}^\tau_{t_0}|\mathbf{y}_{t_0}, \{\mathbf{x}\}^\tau_{t_0}]}{p(\mathbf{y}_\tau)\hat{p}[\mathbf{y}^{t_0}_\tau|\mathbf{y}_\tau, \{\mathbf{x}\}^{t_0}_\tau]}$$
$$+ I_m(\mathbf{x}_{t_0}, \mathbf{y}_{t_0}) + \mathcal{T}_{\mathbf{y}\to\mathbf{x}} - I_m(\mathbf{x}_\tau, \mathbf{y}_\tau) - \mathcal{T}^\dagger_{\mathbf{y}\to\mathbf{x}}$$
$$= \ln\frac{p(\mathbf{x}_{t_0}, \mathbf{y}_{t_0})}{p(\mathbf{x}_\tau, \mathbf{y}_\tau)}\frac{\hat{p}[\mathbf{y}^\tau_{t_0}|\mathbf{y}_{t_0}, \{\mathbf{x}\}^\tau_{t_0}]\hat{p}[\mathbf{x}^\tau_{t_0}|\mathbf{x}_{t_0}, \{\mathbf{y}\}^\tau_{t_0}]}{\hat{p}[\mathbf{y}^{t_0}_\tau|\mathbf{y}_\tau, \{\mathbf{x}\}^{t_0}_\tau]\hat{p}[\mathbf{x}^{t_0}_\tau|\mathbf{x}_\tau, \{\mathbf{y}\}^{t_0}_\tau]}$$
$$- \ln\frac{p(\mathbf{x}_{t_0})}{p(\mathbf{x}_\tau)}\frac{p[\mathbf{x}^\tau_{t_0}|\mathbf{x}_{t_0}]}{p[\mathbf{x}^{t_0}_\tau|\mathbf{x}_\tau]}$$
$$= \ln\frac{p[\mathbf{x}^\tau_{t_0}, \mathbf{y}^\tau_{t_0}]p[\mathbf{x}^{t_0}_\tau]}{p[\mathbf{x}^{t_0}_\tau, \mathbf{y}^{t_0}_\tau]p[\mathbf{x}^\tau_{t_0}]}. \tag{C1}$$

Consequently

$$\left\langle \exp\left[-(\Delta\mathcal{S}^y_{\text{tot}} + \Delta\mathcal{S}^{y|x}_{\text{info}})\right]\right\rangle =$$
$$\int d\mathbf{x}^\tau_{t_0}d\mathbf{y}^\tau_{t_0}p[\mathbf{x}^\tau_{t_0}, \mathbf{y}^\tau_{t_0}]\frac{p[\mathbf{x}^{t_0}_\tau, \mathbf{y}^{t_0}_\tau]p[\mathbf{x}^\tau_{t_0}]}{p[\mathbf{x}^\tau_{t_0}, \mathbf{y}^\tau_{t_0}]p[\mathbf{x}^{t_0}_\tau]}$$
$$= \int d\mathbf{x}^\tau_{t_0}\frac{p[\mathbf{x}^\tau_{t_0}]}{p[\mathbf{x}^{t_0}_\tau]}\int d\mathbf{y}^\tau_{t_0}p[\mathbf{x}^{t_0}_\tau, \mathbf{y}^{t_0}_\tau]$$
$$= \int d\mathbf{x}^\tau_{t_0}\frac{p[\mathbf{x}^\tau_{t_0}]}{p[\mathbf{x}^{t_0}_\tau]}p[\mathbf{x}^{t_0}_\tau]$$
$$= \int d\mathbf{x}^\tau_{t_0}p[\mathbf{x}^\tau_{t_0}] = 1. \tag{C2}$$

The positivity of $\langle\Delta\mathcal{S}^y_{\text{tot}}\rangle + \langle\Delta\mathcal{S}^{y|x}_{\text{info}}\rangle$ then follows from Jensen's inequality.

## Appendix D: Solution to the coupled harmonic oscillators with heat baths at distinct temperatures and $s = 0$ transfer entropy

The probability distribution of the system described by eq. (49) is described by the Fokker Planck equation

$$\frac{\partial p(x,y)}{\partial t} = \frac{\partial}{\partial x}\left(p(x,y)\frac{(k+k_c)x - k_c y}{\gamma_x}\right)$$
$$+ \frac{\partial}{\partial y}\left(p(x,y)\frac{(k+k_c)y - k_c x}{\gamma_y}\right) + \frac{k_B T_x}{\gamma_x}\frac{\partial^2 p(x,y)}{\partial x^2}$$
$$+ \frac{k_B T_y}{\gamma_y}\frac{\partial^2 p(x,y)}{\partial y^2} \tag{D1}$$

which has stationary solution

$$p^s(x,y) = (2\pi|\Psi|)^{-\frac{1}{2}}\exp\left[\mathbf{z}^T\Psi^{-1}\mathbf{z}\right] \tag{D2}$$

where

$$\mathbf{z} = \begin{pmatrix}\mathbf{x}\\\mathbf{y}\end{pmatrix}; \quad \Psi = \begin{pmatrix}\Psi_{xx} & \Psi_{xy}\\\Psi_{yx} & \Psi_{yy}\end{pmatrix} \tag{D3}$$

and

$$\Psi_{xx} = \frac{\gamma_y\left(k_B T_x k(k+2k_c) + k_c^2 k_B T_y\right) + k_B T_x\gamma_x(k+k_c)^2}{k(k+k_c)(k+2k_c)(\gamma_x+\gamma_y)}$$

$$\Psi_{yy} = \frac{\gamma_x\left(k_B T_y k(k+2k_c) + k_c^2 k_B T_x\right) + k_B T_y\gamma_y(k+k_c)^2}{k(k+k_c)(k+2k_c)(\gamma_x+\gamma_y)}$$

$$\Psi_{xy} = \Psi_{yx} = \frac{k_c(k_B T_x\gamma_x + k_B T_y\gamma_y)}{k(k+2k_c)(\gamma_x+\gamma_y)}. \tag{D4}$$

Given such a system we can approximate the transfer entropy rate in the stationary state by starting with the short time propagator [52, 79] which yields

$$p(x_{i+1}|x_i, y_i) = \sqrt{\frac{\gamma_x}{4\pi k_B T_x dt}}$$
$$\times \exp\left[-\frac{\gamma_x}{4k_B T_x dt}\left(dx_i + \frac{(k_x+k_c)x_i - k_c y_i}{\gamma_x}dt\right)^2\right] \tag{D5}$$

where $dx_i = x_{i+1} - x_i$. We can approximate the coarse grained dynamics viz.

$$p^s(x_{i+1}|x_i) = \frac{\int_{-\infty}^{\infty}dy_i p(x_{i+1}|x_i, y_i)p^s(x_i, y_i)}{\int_{-\infty}^{\infty}dy_i p^s(x_i, y_i)}$$
$$= \sqrt{\frac{\gamma_x}{4\pi k_B T dt}}\exp\left[-\frac{\gamma_x}{4k_B T dt}\left(dx_i - \frac{\mathcal{F}(x_i)}{\gamma_x}dt\right)^2\right] + O(dt^2) \tag{D6}$$

where

$$\mathcal{F}(x) = \frac{T_x k(k+k_c)(k+2k_c)(\gamma_x+\gamma_y)x}{(\gamma_y\left(T_x k(k+2k_c) + k_c^2 T_y\right) + T_x\gamma_x(k+k_c)^2)} \tag{D7}$$

The local transfer entropy of such an infinitesimal transition is the logarithm of the ratio of these two quantities which is, up to $O(dt)$ employing the appropriate stochastic calculus such that $dx_i^2 = (2k_B T_x/m\gamma_x)dt$, $x_{i+1}dx_i = x_i dx_i + (2k_B T_x/\gamma_x)dt + O(dt^{3/2})$, $x_{i+1}dt = x_i dt + O(dt^{3/2})$ and $y_{i+1}dx_i = y_i dx_i + O(dt^{3/2})$,

$$t_{y\to x}^{(0)} = \frac{1}{2}\phi^2(x_i,y_i)dt + \phi(x_i,y_i)dW_i^x$$

$$t_{y\to x}^{\dagger,(0)} = \frac{1}{2}\phi^2(x_i,y_i)dt + \phi(x_i,y_i)dW_i^x + \frac{\phi(x_i,y_i)}{k_B T_x}\circ dx_i$$

$$\phi(x_i,y_i) = \sqrt{\frac{1}{2m\gamma_x k_B T_x}}\left(\mathcal{F}(x_i) + (k+k_c)x_i - k_c y_i\right)$$

$$(D8)$$

where $\circ$ indicates Stratonovich integration. A similar approach gives

$$\Delta\mathcal{S}_{\mathrm{med}}^y = \frac{k_c x_i - (k+k_c)y_i}{k_B T_y}\circ dy_i. \qquad (D9)$$

Performing the relevant averages gives

$$\left\langle t_{y\to x}^{(0)}\right\rangle = (4T_x\gamma_x(k+k_c)(\gamma_x+\gamma_y))^{-1}$$
$$\times \frac{k_c^2\left(k_c^2\gamma_x\gamma_y(T_x-T_y)^2 + (k+k_c)^2(\gamma_x+\gamma_y)^2 T_x T_y\right)}{(\gamma_y\left(kT_x(k+2k_c)+k_c^2 T_y\right)+T_x\gamma_x(k+k_c)^2)}$$

$$\left\langle t_{y\to x}^{\dagger,(0)}\right\rangle = (4T_x\gamma_x(k+k_c)(\gamma_x+\gamma_y))^{-1}$$
$$\times \frac{k_c^2\left(k_c^2\gamma_x\gamma_y(T_x-T_y)^2 + (k+k_c)^2(\gamma_x+\gamma_y)^2 T_x T_y\right)}{(\gamma_y\left(kT_x(k+2k_c)+k_c^2 T_y\right)+T_x\gamma_x(k+k_c)^2)}$$
$$+ \frac{k_c^2(T_x-T_y)}{(k+k_c)T_x(\gamma_x+\gamma_y)}. \qquad (D10)$$

[1] U. Seifert, Rep. Prog. Phys. **75**, 126001 (2012).
[2] U. Seifert, Phys. Rev. Lett. **95**, 040602 (2005).
[3] U. Seifert, Eur. Phys. J. B **64**, 423 (2008).
[4] C. Van den Broeck and M. Esposito, Physica A: Statistical Mechanics and its Applications Proceedings of the 13th International Summer School on Fundamental Problems in Statistical Physics, **418**, 6 (2015).
[5] D. J. Evans, E. G. D. Cohen, and G. P. Morriss, Phys. Rev. Lett. **71**, 2401 (1993).
[6] C. Jarzynski, Phys. Rev. E **56**, 5018 (1997).
[7] C. Jarzynski, Phys. Rev. Lett. **78**, 2690 (1997).
[8] G. E. Crooks, Phys. Rev. E **60**, 2721 (1999).
[9] J. L. Lebowitz and H. Spohn, Journal of Statistical Physics **95**, 333 (1999).
[10] R. J. Harris and G. M. Schütz, J. Stat. Mech. **2007**, P07020 (2007).
[11] R. Spinney and I. Ford, in *Nonequilibrium Statistical Physics of Small Systems*, edited by R. Klages, W. Just, and C. Jarzynski (Wiley-VCH Verlag GmbH & Co. KGaA, 2013) pp. 3–56.
[12] J. M. R. Parrondo, C. V. den Broeck, and R. Kawai, New J. Phys. **11**, 073008 (2009).
[13] C. Maes and K. Netočný, Journal of Statistical Physics **110**, 269 (2003).
[14] R. Kawai, J. M. R. Parrondo, and C. Van den Broeck, Phys. Rev. Lett. **98**, 080602 (2007).
[15] P. Gaspard, Journal of Statistical Physics **117**, 599 (2004).
[16] L. Barnett and T. Bossomaier, Phys. Rev. Lett. **109**, 138105 (2012).
[17] J. T. Lizier, M. Prokopenko, and A. Y. Zomaya, Phys. Rev. E **77**, 026110 (2008).
[18] J. T. Lizier, M. Prokopenko, and A. Y. Zomaya, Chaos **20**, 037109 (2010).
[19] J. T. Lizier, M. Prokopenko, and A. Y. Zomaya, Information Sciences **208**, 39 (2012).
[20] J. T. Lizier, in *Directed Information Measures in Neuroscience*, Understanding Complex Systems, edited by M. Wibral, R. Vicente, and J. T. Lizier (Springer Berlin Heidelberg, 2014) pp. 161–193.
[21] M. Lungarella and O. Sporns, PLoS Comput Biol **2** (2006), 10.1371/journal.pcbi.0020144.
[22] X. R. Wang, J. T. Lizier, and M. Prokopenko, Artificial Life **17**, 315 (Fall 2011).
[23] J. M. R. Parrondo, J. M. Horowitz, and T. Sagawa, Nat Phys **11**, 131 (2015).
[24] T. Sagawa and M. Ueda, Phys. Rev. Lett. **104**, 090602 (2010).
[25] S. Toyabe, T. Sagawa, M. Ueda, E. Muneyuki, and M. Sano, Nat Phys **6**, 988 (2010).
[26] D. Abreu and U. Seifert, Phys. Rev. Lett. **108**, 030601 (2012).
[27] T. Sagawa and M. Ueda, Phys. Rev. E **85**, 021104 (2012).
[28] G. Diana and M. Esposito, Journal of Statistical Mechanics: Theory and Experiment **2014**, P04010 (2014).
[29] D. Hartich, A. C. Barato, and U. Seifert, J. Stat. Mech. **2014**, P02016 (2014).
[30] J. M. Horowitz and M. Esposito, Phys. Rev. X **4**, 031015 (2014).
[31] J. M. Horowitz and H. Sandberg, New J. Phys. **16**, 125007 (2014).
[32] D. Mandal and C. Jarzynski, PNAS **109**, 11641 (2012).
[33] J. M. Horowitz, T. Sagawa, and J. M. R. Parrondo, Phys. Rev. Lett. **111**, 010602 (2013).
[34] D. Abreu and U. Seifert, EPL **94**, 10001 (2011).
[35] M. Bauer, D. Abreu, and U. Seifert, J. Phys. A: Math. Theor. **45**, 162001 (2012).
[36] E. Gerstner, Nature News (2002), 10.1038/news020722-2.
[37] G. M. Wang, E. M. Sevick, E. Mittag, D. J. Searles, and D. J. Evans, Phys. Rev. Lett. **89**, 050601 (2002).
[38] T. Schreiber, Phys. Rev. Lett. **85**, 461 (2000).
[39] C. Jarzynski, Annual Review of Condensed Matter Physics **2**, 329 (2011).
[40] J. Schnakenberg, Rev. Mod. Phys. **48**, 571 (1976).
[41] S. Rahav and C. Jarzynski, J. Stat. Mech. **2007**, P09012 (2007).

[42] A. C. Barato, D. Hartich, and U. Seifert, Phys. Rev. E **87**, 042104 (2013).

[43] C. Van den Broeck and M. Esposito, Phys. Rev. E **82**, 011144 (2010).

[44] M. Esposito, Phys. Rev. E **85**, 041125 (2012).

[45] M. Esposito and C. Van den Broeck, Phys. Rev. E **82**, 011143 (2010).

[46] M. Esposito and C. Van den Broeck, Phys. Rev. Lett. **104**, 090601 (2010).

[47] T. Bossomaier, L. Barnett, M. Harré, and J. T. Lizier, *An Introduction to Transfer Entropy: Information Flow in Complex Systems* (Springer, 2016) in press.

[48] L. Barnett and A. K. Seth, "Detectability of Granger causality for subsampled continuous-time neurophysiological processes," (2016), arXiv:1606.08644.

[49] J. L. Lebowitz and P. G. Bergmann, Annals of Physics **1**, 1 (1957).

[50] J. L. Lebowitz, Phys. Rev. **114**, 1192 (1959).

[51] P. G. Bergmann and J. L. Lebowitz, Phys. Rev. **99**, 578 (1955).

[52] R. E. Spinney and I. J. Ford, Phys. Rev. E **85**, 051113 (2012).

[53] R. E. Spinney and I. J. Ford, Phys. Rev. Lett. **108**, 170603 (2012).

[54] S. Bandopadhyay, D. Chaudhuri, and A. M. Jayannavar, Phys. Rev. E **92**, 032143 (2015).

[55] B. Derrida, J. Stat. Mech. **2007**, P07023 (2007).

[56] H. Tasaki, arXiv:0706.1032 [cond-mat] (2007).

[57] C. Maes, Progress of Theoretical Physics Supplement **184**, 318 (2010).

[58] M. Colangeli, C. Maes, and B. Wynants, J. Phys. A: Math. Theor. **44**, 095001 (2011).

[59] M. Bauer and F. Cornu, J. Phys. A: Math. Theor. **48**, 015008 (2015).

[60] J. M. Horowitz, J. Stat. Mech. **2015**, P03006 (2015).

[61] D. Hartich, A. C. Barato, and U. Seifert, Phys. Rev. E **93**, 022116 (2016).

[62] A. Gomez-Marin, J. M. R. Parrondo, and C. Van den Broeck, Phys. Rev. E **78**, 011107 (2008).

[63] K. Kawaguchi and Y. Nakayama, Phys. Rev. E **88**, 022147 (2013).

[64] A. C. Barato, D. Hartich, and U. Seifert, New Journal of Physics **16**, 103024 (2014).

[65] T. M. Cover and J. A. Thomas, *Elements of Information Theory* (John Wiley & Sons, Inc., 2005).

[66] E. H. Feng and G. E. Crooks, Phys. Rev. Lett. **101**, 090602 (2008).

[67] D. Marelli, K. Mahata, and M. Fu, IEEE Transactions on Signal Processing **60**, 3052 (2012).

[68] É. Roldán, I. Neri, M. Dörpinghaus, H. Meyr, and F. Jülicher, Phys. Rev. Lett. **115**, 250602 (2015).

[69] T. B. Batalhão, A. M. Souza, R. S. Sarthour, I. S. Oliveira, M. Paternostro, E. Lutz, and R. M. Serra, Phys. Rev. Lett. **115**, 190601 (2015).

[70] E. H. Feng and G. E. Crooks, Phys. Rev. E **79**, 012104 (2009).

[71] S. Ito and T. Sagawa, Phys. Rev. Lett. **111**, 180603 (2013).

[72] J. Shore and R. Johnson, IEEE Transactions on Information Theory **26**, 26 (1980).

[73] S. Kullback, *Information Theory and Statistics* (John Wiley & Sons, 1959).

[74] G. P. Karev, Entropy **12**, 1673 (2010).

[75] M. Ponmurugan, Phys. Rev. E **82**, 031129 (2010).

[76] M. Prokopenko and J. T. Lizier, Scientific Reports **4** (2014), 10.1038/srep05394.

[77] M. Prokopenko, J. T. Lizier, and D. C. Price, Entropy **15**, 524 (2013).

[78] M. Prokopenko and I. Einav, Phys. Rev. E **91**, 062143 (2015).

[79] C. Wissel, Z Physik B **35**, 185 (1979).