

Article

On Thermodynamic Interpretation of Transfer Entropy

Mikhail Prokopenko^{1,2,*}, Joseph T. Lizier¹ and Don C. Price³

¹ CSIRO Information and Communications Technologies Centre, PO Box 76, Epping, NSW 1710, Australia

² School of Physics, University of Sydney, Australia

³ CSIRO Materials Science and Engineering, Bradfield Road, West Lindfield, NSW 2070, Australia

* Author to whom correspondence should be addressed; mikhail.prokopenko@csiro.au, +61(2) 9372 4716

Version February 6, 2013 submitted to *Entropy*. Typeset by \LaTeX using class file *mdpi.cls*

Abstract: We propose a thermodynamic interpretation of transfer entropy near equilibrium, using a specialised Boltzmann's principle. The approach relates conditional probabilities to the probabilities of the corresponding state transitions. This in turn characterises transfer entropy as a difference of two entropy rates: the rate for a resultant transition and another rate for a possibly irreversible transition within the system affected by an additional source. We then show that this difference, the local transfer entropy, is proportional to the external entropy production, possibly due to irreversibility. Near equilibrium, transfer entropy is also interpreted as the difference in equilibrium stabilities with respect to two scenarios: a default case and the case with an additional source. Finally, we demonstrated that such a thermodynamic treatment is not applicable to information flow, a measure of causal effect.

Keywords: transfer entropy; information transfer; entropy production; irreversibility; Kullback-Leibler divergence; thermodynamic equilibrium; Boltzmann's principle; causal effect

1. Introduction

Transfer entropy has been introduced as an information-theoretic measure that quantifies the statistical coherence between systems evolving in time [1]. Moreover, it was designed to detect asymmetry in the interaction of subsystems by distinguishing between “driving” and “responding” elements. In constructing the measure, Schreiber considered several candidates as measures of directional

19 information transfer, including symmetric mutual information, time-delayed mutual information, as
20 well as asymmetric conditional information. All these alternatives were argued to be inadequate for
21 determining the direction of information transfer between two, possibly coupled, processes.

22 In particular, defining information transfer simply as the dependence of the next state of the receiver
23 on the previous state of the source [2] is incomplete according to Schreiber's criteria requiring the
24 definition to be both *directional* and *dynamic*. Instead, the (predictive) information transfer is defined as
25 the average information contained in the source about the next state of the destination in the context of
26 what was already contained in the destination's past.

27 Following the seminal work of Schreiber [1] numerous applications of transfer entropy have been
28 successfully developed, by capturing information transfer within complex systems, e.g., the stock market
29 [3], food webs [4], EEG signals [5], biochemicals [6], cellular automata and distributed computation in
30 general [7–10], modular robotics [11], random and small-world Boolean networks [12,13], inter-regional
31 interactions within a brain [14], swarm dynamics [15], cascading failures in power grids [16], etc.
32 Also, several studies further capitalised on transition probabilities used in the measure, highlighting
33 fundamental connections of the measure to entropy rate and Kullback-Leibler divergence noted by
34 Kaiser and Schreiber [17], as well as causal flows [18]. At the same time there are several recent
35 studies investigating ties between information theory and thermodynamics [19–23]. This is primarily
36 through Landauer's principle [24], which states that irreversible destruction of one bit of information
37 results in dissipation of at least $kT \ln 2$ J of energy¹ into the environment (i.e. an entropy increase in the
38 environment by this amount).²

39 Nevertheless, transfer entropy *per se* has not been precisely interpreted thermodynamically. Of
40 course, as a measure of directed information transfer, it does not need to have an explicit thermodynamic
41 meaning. Yet, one may still put forward several questions attempting to cast the measure in terms more
42 familiar to a physicist rather than an information theorist or a computer scientist: is transfer entropy a
43 measure of some entropy transferred between subsystems or coupled processes? is it instead an entropy
44 of some transfer happening within the system under consideration? (and what is then the nature of such
45 transfer?). If it is simply a difference between some entropy rates, as can be seen from the definition
46 itself, one may still inquire about the thermodynamic nature of the underlying processes.

47 Obviously, once the subject relating entropy definitions from information theory and thermodynamics
48 is touched, one may expect vigorous debates that have been ongoing since Shannon introduced the
49 term entropy itself. While this paper will attempt to produce a thermodynamic interpretation of transfer
50 entropy, it is out of scope to comment here on rich connections between Boltzmann entropy and Shannon
51 entropy, or provide a review of quite involved discussions on the topic. It suffices to point out prominent
52 works of Jaynes [26,27] who convincingly demonstrated that information theory can be applied to the
53 problem of justification of statistical mechanics, producing predictions of equilibrium thermodynamic
54 properties. The statistical definition of entropy is widely considered more general and fundamental than
55 the original thermodynamic definition, sometimes allowing for extensions to the situations where the
56 system is not in thermal equilibrium [23,28]. In this study, however, we treat the problem of finding a

¹ T is the absolute temperature and k is Boltzmann's constant.

² Maroney [25] argues that while a logically irreversible transformation of information does generate this amount of heat, it can in fact be accomplished by a *thermodynamically reversible* mechanism.

57 thermodynamic interpretation of transfer entropy somewhat separately from the body of work relating
58 Boltzmann and Shannon entropies — and the reason for this is mainly that, even staying within Jaynes’
59 framework, one still needs to provide a possible thermodynamic treatment for transfer entropy *per se*.
60 As will become clear, this task is not trivial, and needs to be approached carefully.

61 Another contribution of this paper is a clarification that similar thermodynamic treatment is not
62 applicable to information flow — a measure introduced by Ay and Polani [18] in order to capture causal
63 effect. That correlation is not causation is well-understood. Yet while authors increasingly consider
64 the notions of information transfer and information flow and how they fit with our understanding of
65 correlation and causality [1,18,29–34], several questions nag. Is information transfer, captured by
66 transfer entropy, akin to causal effect? If not, what is the distinction between them? When examining the
67 “effect” of one variable on another (e.g. between brain regions), should one seek to measure information
68 transfer or causal effect?

69 Unfortunately, these concepts have become somewhat tangled in discussions of information transfer.
70 Measures for both predictive transfer [1] and causal effect [18] have been inferred to capture information
71 transfer in general, and measures of predictive transfer have been used to infer causality [33,35–37]
72 with the two sometimes (problematically) directly equated (e.g. [29,32,34,38–40]). The study of Lizier
73 and Prokopenko [41] clarified the relationship between these concepts and described the manner in
74 which they should be considered separately. Here, in addition, we demonstrate that a thermodynamic
75 interpretation of transfer entropy is not applicable to causal effect (information flow), and clarify the
76 reasons behind this.

77 This paper is organised as follows. We begin with Section 2 that introduces relevant information-
78 theoretic measures both in average and local terms. Section 3 defines the system and the range of
79 applicability of our approach. In providing a thermodynamic interpretation for transfer entropy in
80 Section 4 we relate conditional probabilities to the probabilities of the corresponding state transitions,
81 and use a specialised Boltzmann’s principle. This allows us to define components of transfer entropy
82 with the entropy rate of (i) the resultant transition and (ii) the internal entropy production. Sub-section
83 4.3 presents an interpretation of transfer entropy near equilibrium. The following Section 5 discusses the
84 challenges for supplying a similar interpretation to causal effect (information flow). A brief discussion
85 in Section 6 concludes the paper.

86 2. Definitions

87 In the following sections we describe relevant background on transfer entropy and causal effect
88 (information flow), along some technical preliminaries.

89 2.1. Transfer entropy

90 Mutual information $I_{Y;X}$ has been something of a de facto measure for information transfer between
91 Y and X in complex systems science in the past (e.g. [42–44]). A major problem however is that mutual
92 information contains no inherent *directionality*. Attempts to address this include using the previous state
93 of the “source” variable Y and the next state of the “destination” variable X' (known as *time-lagged*

94 *mutual information* $I_{Y;X'}$). However, Schreiber [1] points out that this ignores the more fundamental
 95 problem that mutual information measures the *statically* shared information between the two elements.³

To address these inadequacies Schreiber introduced *transfer entropy* [1] (TE), the deviation from independence (in bits) of the state transition (from the previous state to the next state) of an information destination X from the previous state of an information source Y :

$$T_{Y \rightarrow X}(k, l) = \sum_{x_{n+1}, x_n^{(k)}, y_n^{(l)}} p(x_{n+1}, x_n^{(k)}, y_n^{(l)}) \log_2 \frac{p(x_{n+1} | x_n^{(k)}, y_n^{(l)})}{p(x_{n+1} | x_n^{(k)})}. \quad (1)$$

Here n is a time index, $x_n^{(k)}$ and $y_n^{(l)}$ represent past states of X and Y (i.e. the k and l past values of X and Y up to and including time n). Schreiber points out that this formulation is a truly *directional, dynamic* measure of information transfer, and is a generalisation of the entropy rate to more than one element to form a mutual information *rate*. That is, transfer entropy may be seen as the difference between two entropy rates:

$$T_{Y \rightarrow X}(k, l) = h_X - h_{X,Y}, \quad (2)$$

where h_X is the entropy rate:

$$h_X = - \sum p(x_{n+1}, x_n^{(k)}) \log_2 p(x_{n+1} | x_n^{(k)}), \quad (3)$$

and $h_{X,Y}$ is a generalised entropy rate conditioning on the source state as well:

$$h_{X,Y} = - \sum p(x_{n+1}, x_n^{(k)}, y_n^{(l)}) \log_2 p(x_{n+1} | x_n^{(k)}, y_n^{(l)}). \quad (4)$$

The entropy rate h_X accounts for the average number of bits needed to encode one additional state of the system if all previous states are known [1], while the entropy rate $h_{X,Y}$ is the entropy rate capturing the average number of bits required to represent the value of the next destination's state if source states are included in addition. Since one can always write

$$h_X = - \sum p(x_{n+1}, x_n^{(k)}) \log_2 p(x_{n+1} | x_n^{(k)}) = - \sum p(x_{n+1}, x_n^{(k)}, y_n^{(l)}) \log_2 p(x_{n+1} | x_n^{(k)}), \quad (5)$$

it is easy to see that the entropy rate h_X is equivalent to the rate $h_{X,Y}$ when the next state of destination is independent of the source [1]:

$$p(x_{n+1} | x_n^{(k)}) = p(x_{n+1} | x_n^{(k)}, y_n^{(l)}), \quad (6)$$

96 Thus, in this case the transfer entropy reduces to zero.

Similarly, the TE can be viewed as a *conditional* mutual information $I(Y^{(l)}; X' | X^{(k)})$ [17], that is as the average information contained in the source about the next state X' of the destination that was not already contained in the destination's past $X^{(k)}$:

$$T_{Y \rightarrow X}(k, l) = I_{Y^{(l)}; X' | X^{(k)}} = H_{X' | X^{(k)}} - H_{X' | X^{(k)}, Y^{(l)}}. \quad (7)$$

97 This could be interpreted (following [45] and [44]) as the diversity of state transitions in the destination
 98 minus assortative noise between those state transitions and the state of the source.

³ The same criticism applies to equivalent non information-theoretic definitions such as that in [2].

99 Furthermore, we note that Schreiber’s original description can be rephrased as the information
 100 provided by the source about the state transition in the destination. That $x_n^{(k)} \rightarrow x_{n+1}$ (or including
 101 redundant information $x_n^{(k)} \rightarrow x_{n+1}^{(k)}$) is a *state transition* is underlined in that the $x_n^{(k)}$ are *embedding*
 102 *vectors* [46], which capture the underlying *state* of the process. Indeed, since all of the above
 103 mathematics for the transfer entropy is equivalent if we consider the next source *state* $x_{n+1}^{(k)}$ instead
 104 of the next source value x_{n+1} , we shall adjust our notation from here onwards to consider the next source
 105 state $x_{n+1}^{(k)}$, so that we are always speaking about interactions between source states \mathbf{y}_n and destination
 106 state transitions $\mathbf{x}_n \rightarrow \mathbf{x}_{n+1}$ (with embedding lengths l and k implied).

107 Importantly, the TE remains a measure of observed (conditional) *correlation* rather than direct effect.
 108 In fact, the TE is a non-linear extension of a concept known as the “Granger causality” [47], the
 109 nomenclature for which may have added to the confusion associating information transfer and causal
 110 effect. Importantly, as an information-theoretic measure based on observational probabilities, the TE is
 111 applicable to both deterministic and stochastic systems.

112 2.2. Local transfer entropy

113 Information-theoretic variables are generally defined and used as an *average* uncertainty or
 114 information. We are interested in considering *local* information-theoretic values, i.e. the uncertainty
 115 or information associated with a *particular observation* of the variables rather than the average over all
 116 observations. *Local* information-theoretic measures are sometimes called *point-wise* measures [48,49].
 117 Local measures within a global average are known to provide important insights into the *dynamics* of
 118 non-linear systems [50].

Using the technique originally described in [7], we observe that the TE is an average (or *expectation value*) of a *local transfer entropy* at each observation n , i.e.:

$$T_{Y \rightarrow X} = \langle t_{Y \rightarrow X}(n+1) \rangle, \quad (8)$$

$$t_{Y \rightarrow X}(n+1) = \log_2 \frac{p(\mathbf{x}_{n+1} | \mathbf{x}_n, \mathbf{y}_n)}{p(\mathbf{x}_{n+1} | \mathbf{x}_n)}, \quad (9)$$

119 with embedding lengths l and k implied as described above. The local transfer entropy quantifies the
 120 information contained in the source state \mathbf{y}_n about the next state of the destination \mathbf{x}_{n+1} at time step
 121 $n+1$, in the context of what was already contained in the past state of the destination \mathbf{x}_n . The measure is
 122 *local* in that it is defined at each time n for each destination X in the system and each causal information
 123 source Y of the destination.

The local TE may also be expressed as a local conditional mutual information, or a difference between local conditional entropies:

$$t_{Y \rightarrow X}(n+1) = i(\mathbf{y}_n; \mathbf{x}_{n+1} | \mathbf{x}_n) = h(\mathbf{x}_{n+1} | \mathbf{x}_n) - h(\mathbf{x}_{n+1} | \mathbf{x}_n, \mathbf{y}_n). \quad (10)$$

where local conditional mutual information is given by

$$i(\mathbf{y}_n; \mathbf{x}_{n+1} | \mathbf{x}_n) = \log_2 \frac{p(\mathbf{x}_{n+1} | \mathbf{x}_n, \mathbf{y}_n)}{p(\mathbf{x}_{n+1} | \mathbf{x}_n)} \quad (11)$$

and local conditional entropies are defined analogously:

$$h(\mathbf{x}_{n+1} | \mathbf{x}_n) = -\log_2 p(\mathbf{x}_{n+1} | \mathbf{x}_n), \quad (12)$$

$$h(\mathbf{x}_{n+1} \mid \mathbf{x}_n, \mathbf{y}_n) = -\log_2 p(\mathbf{x}_{n+1} \mid \mathbf{x}_n, \mathbf{y}_n), \quad (13)$$

124 The average transfer entropy $T_{Y \rightarrow X}(k)$ is always positive but is bounded above by the information
 125 capacity of a single observation of the destination. For a discrete system with b possible observations
 126 this is $\log_2 b$ bits. As a conditional mutual information, it can be either larger *or* smaller than the
 127 corresponding mutual information [51]. The *local* TE however is not constrained so long as it averages
 128 into this range: it can be greater than $\log_2 b$ for a large local information transfer, and can also in fact
 129 be measured to be negative. Local transfer entropy is negative where (in the context of the history of
 130 the destination) the probability of observing the actual next state of the destination given the source
 131 state $p(\mathbf{x}_{n+1} \mid \mathbf{x}_n, \mathbf{y}_n)$, is lower than that of observing that actual next state independently of the source
 132 $p(\mathbf{x}_{n+1} \mid \mathbf{x}_n)$. In this case, the source variable is actually *misinformative* or misleading about the state
 133 transition of the destination. It is possible for the source to be misleading where other causal information
 134 sources influence the destination, or in a stochastic system. Full examples are described by Lizier et al.
 135 [7].

136 2.3. Causal effect as information flow

137 As noted earlier, predictive information transfer refers to the amount of information that a source
 138 variable adds to the next state of a destination variable; i.e. “if I know the state of the source, how much
 139 does that help to predict the state of the destination?”. Causal effect, on the contrary, refers to the extent
 140 to which the source variable has a direct influence or drive on the next state of a destination variable,
 141 i.e. “if I change the state of the source, to what extent does that alter the state of the destination?”.
 142 Information from causal effect can be seen to *flow* through the system, like injecting dye into a river
 143 [18].

144 It is well-recognised that measurement of causal effect necessitates some type of *perturbation* or
 145 *intervention* of the source so as to detect the effect of the intervention on the destination (e.g. see
 146 [52]). Attempting to infer causality without doing so leaves one measuring correlations of observations,
 147 regardless of how directional they may be [18]. In this section, we adopt the measure information flow
 148 for this purpose, and describe a method introduced by Lizier and Prokopenko [41] for applying it on a
 149 local scale.

150 Following Pearl’s probabilistic formulation of causal Bayesian networks [52], Ay and Polani [18]
 151 consider how to measure causal information flow via *interventional conditional probability distribution*
 152 *functions*. For instance, an interventional conditional PDF $p(y \mid \hat{s})$ considers the distribution of y
 153 resulting from *imposing* the value of \hat{s} . *Imposing* means intervening in the system to *set* the value of the
 154 imposed variable, and is at the essence of the definition of causal information flow. As an illustration of
 155 the difference between interventional and standard conditional PDFs, consider two correlated variables
 156 S and Y : their correlation alters $p(y \mid s)$ in general from $p(y)$. If both variables are solely caused by
 157 another variable G however, then even where they remain correlated we have $p(y \mid \hat{s}) = p(y)$ because
 158 imposing a value \hat{s} has no effect on the value of y .

159 In a similar fashion to the definition of transfer entropy as the deviation of a destination from
 160 *stochastic* independence on the source in the content of the destination’s past, Ay and Polani propose

161 the measure *information flow* as the deviation of the destination X from *causal* independence on the
 162 source Y *imposing* another set of nodes \mathbf{S} . Mathematically, this is written as:

$$I_p(Y \rightarrow X | \hat{\mathbf{S}}) = \sum_{\mathbf{s}} p(\mathbf{s}) \sum_y p(y | \hat{\mathbf{s}}) \sum_x p(x | \hat{y}, \hat{\mathbf{s}}) \log_2 \frac{p(x | \hat{y}, \hat{\mathbf{s}})}{\sum_{y'} p(y' | \hat{\mathbf{s}}) p(x | \hat{y}', \hat{\mathbf{s}})}. \quad (14)$$

163 The value of the measure is dependent on the choice of the set of nodes \mathbf{S} . It is possible to obtain
 164 a measure of apparent causal information flow $I_p(Y \rightarrow X)$ from Y to X without any \mathbf{S} (i.e. $\mathbf{S} = \emptyset$),
 165 yet this can be misleading. In particular, it ignores causal information flow arising from interactions
 166 of the source with another source variable. For example, if $x = y \text{ XOR } s$ and $p(y, s) = 0.25$ for
 167 each combination of binary y and s , then $I_p(Y \rightarrow X) = 0$ despite the clear causal effect of Y , while
 168 $I_p(Y \rightarrow X | \hat{S}) = 1$ bit. Also, we may have $I_p(Y \rightarrow X) > 0$ only because Y effects \mathbf{S} which in turn
 169 effects X ; where we are interested in *direct* causal information flow from Y to X only $I_p(Y \rightarrow X | \hat{\mathbf{S}})$
 170 validly infers no direct causal effect.

171 Here we are interested in measuring the *direct* causal information flow from Y to X , so we must either
 172 include all possible other sources in \mathbf{S} or at least include enough sources to “block”⁴ all non-immediate
 173 directed paths from Y to X [18]. The minimum to satisfy this is the set of all direct causal sources of
 174 X excluding Y , including any past states of X that are direct causal sources. That is, in alignment with
 175 transfer entropy \mathbf{S} would include $X^{(k)}$.

176 The major task in computing $I_p(Y \rightarrow X | \hat{\mathbf{S}})$ is the determination of the underlying interventional
 177 conditional PDFs in Eq. (14). By definition these may be gleaned by observing the results of intervening
 178 in the system, however this is not possible in many cases.

179 One alternative is to use detailed knowledge of the dynamics, in particular the structure of the causal
 180 links and possibly the underlying rules of the causal interactions. This also is often not available in
 181 many cases, and indeed is often the very goal for which one turned to such analysis in the first place.
 182 Regardless, where such knowledge is available it may allow one to make direct inferences.

183 Under certain constrained circumstances, one can construct these values from observational
 184 probabilities only [18], e.g. with the “back-door adjustment” [52]. A particularly important constraint
 185 on using the back-door adjustment here is that *all* $\{s, y\}$ combinations must be observed.

186 2.4. Local information flow

187 A *local information flow* can be defined following the argument that was used to define local
 188 information transfer:

$$f(y \rightarrow x | \hat{\mathbf{s}}) = \log_2 \frac{p(x | \hat{y}, \hat{\mathbf{s}})}{\sum_{y'} p(y' | \hat{\mathbf{s}}) p(x | \hat{y}', \hat{\mathbf{s}})}. \quad (15)$$

189 The meaning of the local information flow is slightly different however. Certainly, it is an *attribution*
 190 of local causal effect of y on x were $\hat{\mathbf{s}}$ imposed at the given observation (y, x, \mathbf{s}) . However, one must

⁴ A set of nodes U *blocks* a path of causal links where there is a node v on the path such that either:

1. $v \in U$ and the causal links through v on the path are not both into v , or
2. the causal links through v on the path are both into v , and v and all its causal descendants are not in U .

191 be aware that $I_p(Y \rightarrow X | \hat{S})$ is not the *average* of the local values $f(y \rightarrow x | \hat{s})$ in exactly the same
 192 manner as the local values derived for information transfer. Unlike standard information-theoretical
 193 measures, the information flow is averaged over a product of *interventional* conditional probabilities
 194 $(p(s)p(y | \hat{s})p(x | \hat{y}, \hat{s}))$, see Eq. (14) which in general does not reduce down to the probability of the
 195 given observation $p(s, y, x) = p(s)p(y | s)p(x | y, s)$. For instance, it is possible that not all of the
 196 tuples $\{y, x, s\}$ will actually be observed, so averaging over observations would ignore the important
 197 contribution that any unobserved tuples provide to the determination of information flow. Again, the
 198 local information flow is specifically tied not to the *given observation* at time step n but to the *general*
 199 *configuration* (y, x, s) , and only *attributed* to the associated observation of this configuration at time n .

200 3. Preliminaries

201 3.1. System definition

202 Let us consider the non-equilibrium thermodynamics of a physical system close to equilibrium. At
 203 any given moment in time, n , the thermodynamic state of the physical system X is given by a vector
 204 $\mathbf{x} \in R^d$, comprising d variables, for instance the (local) pressure, temperature, chemical concentrations
 205 and so on. A state vector completely describes the physical macrostate as far as predictions of the
 206 outcomes of all possible measurements performed on the system are concerned [53]. The state space of
 207 the system is the set of all possible states of the system.

208 The thermodynamic state is generally considered as a fluctuating entity so that transition probabilities
 209 like $p(\mathbf{x}_{n+1} | \mathbf{x}_n)$ are clearly defined and can be related to a sampling procedure. Each macrostate can
 210 be realised by a number of different microstates consistent with the given thermodynamic variables.
 211 Importantly, in the theory of non-equilibrium thermodynamics close to equilibrium, the microstates
 212 belonging to one macrostate \mathbf{x} are equally probable.

213 3.2. Entropy definitions

The thermodynamic entropy was originally defined by Clausius' as a state function S which satisfies

$$S_B - S_A = \int_A^B dq_{rev}/T, \quad (16)$$

214 where q_{rev} is the heat transferred to an equilibrium thermodynamic system during a reversible process
 215 from state A to state B . Note that this *path integral* is the same for all reversible paths between the past
 216 and next states.

217 It was shown by Jaynes that thermodynamic entropy could be interpreted, from the perspective of
 218 statistical mechanics, as a measure of the amount of information about the microstate of a system that an
 219 observer lacks if they know only the macrostate of the system [53].

220 This is encapsulated in the famous Boltzmann's equation $S = k \log W$, where k is Boltzmann's
 221 constant and W is the number of microstates corresponding to a given macrostate (an integer greater
 222 than or equal to one). While it is not a mathematical probability between zero and one, it is sometimes

223 called “thermodynamic probability”, noting that W can be normalized to a probability $p = W/N$, where
 224 N is the number of possible microstates for all macrostates.

225 The Shannon entropy that corresponds to the Boltzmann entropy $S = k \log W$ is the uncertainty in
 226 the microstate which has produced the given macrostate. That is, given the number W of microscopic
 227 configurations that correspond to the given macrostate, we have $p_i = 1/W$ for each equiprobable
 228 microstate i . As such, we can compute the local entropy for each of these W microstates as
 229 $-\log_2 1/W = \log_2 W$ bits. Note that the average entropy across all of these equiprobable microstates is
 230 $\log_2 W$ bits also. This is equivalent to the Boltzmann entropy up to Boltzmann’s constant k and the base
 231 of the logarithms (see [54,55] for more details).

232 3.3. Transition probabilities

A specialisation of Boltzmann’s principle by Einstein [56], for two states with entropies S and S_0 and
 “relative probability” W_r (the ratio of numbers W and W_0 that account for the numbers of microstates
 in the macrostates with S and S_0 respectively), is given by:

$$S - S_0 = k \log W_r, \quad (17)$$

233 The expression in these relative terms is important, as pointed out by Norton [57], because the probability
 234 W_r is the probability of the transition between the two states under the system’s normal time evolution.

In the example considered by Einstein [56,57], S_0 is the entropy of an (equilibrium) state, e.g. “a
 volume V_0 of space containing n non-interacting, moving points, whose dynamics are such as to favor
 no portion of the space over any other”, while S is the entropy of the (non-equilibrium) state with the
 “same system of points, but now confined to a sub-volume V of V_0 ”. Specifically, Einstein defined the
 transition probability $W_r = (V/V_0)^n$, yielding

$$S - S_0 = kn \log(V/V_0). \quad (18)$$

235 Since dynamics favour no portion of the space over any other, all the microstates are equiprobable.

236 3.4. Entropy production

In general, the variation of entropy of a system ΔS is equal to the sum of the internal entropy
 production σ inside the system and the entropy change due to the interactions with the surroundings
 ΔS_{ext} :

$$\Delta S = \sigma + \Delta S_{ext}, \quad (19)$$

In the case of a closed system, ΔS_{ext} is given by the expression

$$\Delta S_{ext} = \int dq/T, \quad (20)$$

where q represents the heat flow received by the system from the exterior and T is the temperature of the
 system. This expression is often written as

$$\sigma = \Delta S - \Delta S_{ext} = (S - S_0) - \Delta S_{ext}, \quad (21)$$

so that when the transition from the initial state S_0 to the final state S is irreversible, the entropy production $\sigma > 0$, while for reversible processes $\sigma = 0$, that is

$$S - S_0 = \int dq_{rev}/T. \quad (22)$$

237 We shall consider another state vector, \mathbf{y} , describing a state of a part Y of the exterior possibly coupled
 238 to the system represented by X . In other words, X and Y may or may not be dependent. In general, we
 239 shall say that σ_y is the internal entropy production *in the context of* some source Y , while ΔS_{ext} is the
 240 entropy production *attributed to* Y .

241 Alternatively, one may consider two scenarios for such a general physical system. In the first scenario,
 242 the entropy changes only due to reversible transitions, amounting to $S - S_0$. In the second scenario, the
 243 entropy changes partly irreversibly due to the interactions with the external environment affected by \mathbf{y} ,
 244 but still achieves the same total change $S - S_0$.

245 3.5. Range of applicability

In an attempt to provide a thermodynamic interpretation of transfer entropy we make two important assumptions, defining the range of applicability for such an interpretation. The first one relates the transition probability W_{r_1} of the system's reversible state change to the conditional probability $p(\mathbf{x}_{n+1} | \mathbf{x}_n)$, obtained by sampling the process X :

$$p(\mathbf{x}_{n+1} | \mathbf{x}_n) = \frac{1}{Z_1} W_{r_1}, \quad (23)$$

where Z_1 is a normalisation factor which depends on \mathbf{x}_n . According to the expression for transition probability (17), under this assumption the conditional probability of the system's transition from state \mathbf{x}_n to state \mathbf{x}_{n+1} corresponds to some number W_{r_1} , such that $S(\mathbf{x}_{n+1}) - S(\mathbf{x}_n) = k \log W_{r_1}$, and hence

$$p(\mathbf{x}_{n+1} | \mathbf{x}_n) = \frac{1}{Z_1} e^{(S(\mathbf{x}_{n+1}) - S(\mathbf{x}_n))/k}. \quad (24)$$

The second assumption relates the transition probability W_{r_2} of the system's possibly irreversible internal state change, due to the interactions with the external surroundings represented in the state vector \mathbf{y} , to the conditional probability $p(\mathbf{x}_{n+1} | \mathbf{x}_n, \mathbf{y}_n)$, obtained by sampling the systems X and Y :

$$p(\mathbf{x}_{n+1} | \mathbf{x}_n, \mathbf{y}_n) = \frac{1}{Z_2} W_{r_2}. \quad (25)$$

Under this assumption the conditional probability of the system's (irreversible) transition from state \mathbf{x}_n to state \mathbf{x}_{n+1} in the context of \mathbf{y}_n , corresponds to some number W_{r_2} , such that $\sigma_y = k \log W_{r_2}$, where σ_y is the system's internal entropy production in the context of \mathbf{y} , and thus

$$p(\mathbf{x}_{n+1} | \mathbf{x}_n, \mathbf{y}_n) = \frac{1}{Z_2} e^{\sigma_y/k}. \quad (26)$$

246 where Z_2 is a normalisation factor which depends on \mathbf{x}_n .

247 **3.6. An example: random fluctuation near equilibrium**

Let us consider the above-defined stochastic process X for a small random fluctuation around equilibrium:

$$\mathbf{x}_{n+1} = \Lambda \mathbf{x}_n + \xi, \quad (27)$$

248 where ξ is a multi-variate Gaussian noise process, with covariance matrix Σ_ξ , uncorrelated in time.
249 Starting at time n with state \mathbf{x}_n having entropy $S(\mathbf{x}_n)$, the state develops into \mathbf{x}_{n+1} , with entropy
250 $S(\mathbf{x}_{n+1})$.

From the probability distribution function of the above multi-variate Gaussian process, we obtain

$$p(\mathbf{x}_{n+1} | \mathbf{x}_n) = \frac{1}{Z} e^{-\frac{1}{2}(\mathbf{x}_{n+1} - \Lambda \mathbf{x}_n)^T \Sigma_\xi^{-1} (\mathbf{x}_{n+1} - \Lambda \mathbf{x}_n)}. \quad (28)$$

We now demonstrate that this expression concurs with the corresponding expression obtained under assumption (24). To do so we expand the entropies around $\mathbf{x} = 0$ with entropy $S(0)$:

$$S(\mathbf{x}_n) = S(0) - k \frac{1}{2} \mathbf{x}_n^T \Sigma_x^{-1} \mathbf{x}_n \quad (29)$$

251 where Σ_x is the covariance matrix of the process X .

Then, according to the assumption (24)

$$p(\mathbf{x}_{n+1} | \mathbf{x}_n) = \frac{1}{Z_1} e^{(S(\mathbf{x}_{n+1}) - S(\mathbf{x}_n))/k} = \frac{1}{Z_1} e^{-\frac{1}{2}(\mathbf{x}_{n+1}^T \Sigma_x^{-1} \mathbf{x}_{n+1} - \mathbf{x}_n^T \Sigma_x^{-1} \mathbf{x}_n)} = \frac{1}{\tilde{Z}_1} e^{-\frac{1}{2} \mathbf{x}_{n+1}^T \Sigma_x^{-1} \mathbf{x}_{n+1}}, \quad (30)$$

where the term $e^{\frac{1}{2} \mathbf{x}_n^T \Sigma_x^{-1} \mathbf{x}_n}$ is absorbed into the normalisation factor being only dependent on \mathbf{x}_n . In general [58,59], we have

$$\Sigma_x = \sum_{j=0}^{\infty} \Lambda^j \Sigma_\xi \Lambda^{jT}. \quad (31)$$

Given the quasistationarity of the relaxation process, assumed near an equilibrium, $\Lambda \rightarrow 0$, and hence $\Sigma_x \rightarrow \Sigma_\xi$. Then the equation (30) reduces to

$$p(\mathbf{x}_{n+1} | \mathbf{x}_n) = \frac{1}{\tilde{Z}_1} e^{-\frac{1}{2}(\mathbf{x}_{n+1}^T \Sigma_\xi^{-1} \mathbf{x}_{n+1})}. \quad (32)$$

252 The last expression concurs with (28) when $\Lambda \rightarrow 0$.

253 **4. Transfer entropy: thermodynamic interpretation**254 **4.1. Transitions near equilibrium**

255 Supported by this background, we proceed to interpret transfer entropy via transitions between states.
256 In doing so, we shall operate with local information theoretic measures (such as the local transfer entropy
257 (9)), as we are dealing with (transitions between) *specific* states \mathbf{y}_n , \mathbf{x}_n , \mathbf{x}_{n+1} , etc. and not with all
258 possible state-spaces X , Y , etc. containing all realizations of specific states.

Transfer entropy is a difference not between entropies, but rather between entropy rates or conditional entropies, specified on average by (2) or (7), or for local values by (10):

$$t_{Y \rightarrow X}(n+1) = h(\mathbf{x}_{n+1} | \mathbf{x}_n) - h(\mathbf{x}_{n+1} | \mathbf{x}_n, \mathbf{y}_n). \quad (33)$$

As mentioned above, the first assumption (23), taken to define the range of applicability for our interpretation, entails (24). It then follows that the first component of equation (33), $h(\mathbf{x}_{n+1} | \mathbf{x}_n)$, accounts for $S(\mathbf{x}_{n+1}) - S(\mathbf{x}_n)$:

$$h(\mathbf{x}_{n+1} | \mathbf{x}_n) = -\log_2 p(\mathbf{x}_{n+1} | \mathbf{x}_n) = -\log_2 \frac{1}{Z_1} e^{(S(\mathbf{x}_{n+1}) - S(\mathbf{x}_n))/k} \quad (34)$$

$$= \log_2 Z_1 - \frac{1}{k \log 2} (S(\mathbf{x}_{n+1}) - S(\mathbf{x}_n)). \quad (35)$$

259 That is, the local conditional entropy $h(\mathbf{x}_{n+1} | \mathbf{x}_n)$ corresponds to resultant entropy change of the
260 transition from the past state \mathbf{x}_n to the next state \mathbf{x}_{n+1} .

261 Now we need to interpret the second component of (33): the local conditional entropy $h(\mathbf{x}_{n+1} |$
262 $\mathbf{x}_n, \mathbf{y}_n)$ in presence of some other factor or extra source, \mathbf{y}_n . Importantly, we must keep both the past
263 state \mathbf{x}_n and the next state \mathbf{x}_{n+1} the same — only then we can characterise the internal entropy change,
264 offset by some contribution of the source \mathbf{y}_n .

Our second constraint on the system (25) entails (26), and so

$$h(\mathbf{x}_{n+1} | \mathbf{x}_n, \mathbf{y}_n) = -\log_2 p(\mathbf{x}_{n+1} | \mathbf{x}_n, \mathbf{y}_n) = -\log_2 \frac{1}{Z_2} e^{\sigma_y/k} = \log_2 Z_2 - \frac{1}{k \log 2} (\sigma_y). \quad (36)$$

265 4.2. Transfer entropy as entropy production

At this stage we can bring two right-hand side components of transfer entropy (33), represented by (35) and (36), together:

$$t_{Y \rightarrow X}(n+1) = \log_2 \frac{Z_1}{Z_2} + \frac{1}{k \log 2} (-(S(\mathbf{x}_{n+1}) - S(\mathbf{x}_n)) + \sigma_y). \quad (37)$$

When one considers a small fluctuation near an equilibrium, $Z_1 \approx Z_2$, as the number of microstates does not change much in the relevant macrostates. This removes the additive constant. Then, using the expression for entropy production (21), we obtain

$$t_{Y \rightarrow X}(n+1) = -\frac{\Delta S_{ext}}{k \log 2}. \quad (38)$$

266 If $Z_1 \neq Z_2$, the relationship includes some additive constant $\log_2 \frac{Z_1}{Z_2}$.

267 That is, the transfer entropy is proportional to the external entropy production, brought about by the
268 source of irreversibility Y . It captures the difference between the entropy rates that correspond to two
269 scenarios: the reversible process and the irreversible process affected by another source Y . It is neither
270 a transfer of entropy, nor an entropy of some transfer — it is formally a difference between two entropy
271 rates. The opposite sign reflects the different direction of entropy production attributed to the source Y :
272 when $\Delta S_{ext} > 0$, i.e. the entropy increased during the transition in X more than the entropy produced
273 internally, then the local transfer entropy is negative, and the source misinforms about the macroscopic

274 state transition. When, on the other hand, $\Delta S_{ext} < 0$, i.e. some of the internal entropy produced
 275 during the transition in X dissipated to the exterior, then the local transfer entropy is positive, and better
 276 predictions can be made about the macroscopic state transitions in X if source Y is measured.

277 As mentioned earlier, while transfer entropy is non-negative on average, some local transfer
 278 entropies can be negative when (in the context of the history of the destination) the source variable
 279 is misinformative or misleading about the state transition. This, obviously, concurs with the fact that,
 280 while a statistical ensemble average of time averages of the entropy change is always non-negative, at
 281 certain times entropy change can be negative. This follows from the fluctuation theorem [60], the Second
 282 law inequality [61], and can be illustrated with other examples of backward transformations and local
 283 violations of the second law [62,63].

284 Another observation follows from our assumptions (24), (26) and the representation (37) when $Z_1 \approx$
 285 Z_2 . If the local conditional entropy $h(\mathbf{x}_{n+1} | \mathbf{x}_n)$, corresponding to the resultant entropy change of the
 286 transition, is different from the local conditional entropy $h(\mathbf{x}_{n+1} | \mathbf{x}_n, \mathbf{y}_n)$ capturing the internal entropy
 287 production in context of the external source Y , then X and Y are dependent. Conversely, whenever these
 288 two conditional entropies are equal to each other, X and Y are independent.

289 4.3. Transfer entropy as a measure of equilibrium's stability

There is another possible interpretation that considers a fluctuation near the equilibrium. Using
 Kullback-Leibler divergence between discrete probability distributions p and q :

$$D_{\text{KL}}(p||q) = \sum_i p(i) \log \frac{p(i)}{q(i)}, \quad (39)$$

and its local counterpart:

$$d_{\text{KL}}(p||q) = \log \frac{p(i)}{q(i)}, \quad (40)$$

we may also express the local conditional entropy as follows:

$$h(\mathbf{x}_{n+1} | \mathbf{x}_n) = h(\mathbf{x}_{n+1}, \mathbf{x}_n) - h(\mathbf{x}_n) = d_{\text{KL}}(p(\mathbf{x}_{n+1}, \mathbf{x}_n)||p(\mathbf{x}_n)). \quad (41)$$

290 It is known in macroscopic thermodynamics that stability of an equilibrium can be measured with
 291 Kullback-Leibler divergence between the initial (past) state, e.g. \mathbf{x}_n , and the state brought about by
 292 some fluctuation (a new observation), e.g. \mathbf{x}_{n+1} [64]. That is, we can also interpret the local conditional
 293 entropy $h(\mathbf{x}_{n+1} | \mathbf{x}_n)$ as the entropy change (or entropy rate) of the fluctuation near the equilibrium.

Analogously, the entropy change in another scenario, where an additional source \mathbf{y} contributes to the
 fluctuation around the equilibrium, corresponds now to Kullback-Leibler divergence

$$h(\mathbf{x}_{n+1} | \mathbf{x}_n, \mathbf{y}_n) = h(\mathbf{x}_{n+1}, \mathbf{x}_n, \mathbf{y}_n) - h(\mathbf{x}_n, \mathbf{y}_n) = d_{\text{KL}}(p(\mathbf{x}_{n+1}, \mathbf{x}_n, \mathbf{y}_n)||p(\mathbf{x}_n, \mathbf{y}_n)) \quad (42)$$

294 and can be seen as a measure of stability with respect to the fluctuation that is now affected by the extra
 295 source \mathbf{y} .

Contrasting both these fluctuations around the same equilibrium, we obtain in terms of Kullback-
 Leibler divergences:

$$t_{Y \rightarrow X}(n+1) = d_{\text{KL}}(p(\mathbf{x}_{n+1}, \mathbf{x}_n)||p(\mathbf{x}_n)) - d_{\text{KL}}(p(\mathbf{x}_{n+1}, \mathbf{x}_n, \mathbf{y}_n)||p(\mathbf{x}_n, \mathbf{y}_n)). \quad (43)$$

296 In these terms, transfer entropy contrasts stability of the equilibrium between two scenarios: the first
 297 one corresponds to the original system, and the second one disturbs the system by the source Y . If, for
 298 instance, the source Y is such that the system X is independent of it, then there is no difference in the
 299 extents of disturbances to the equilibrium, and the transfer entropy is zero.

300 4.4. Heat transfer

301 It is possible to provide a similar thermodynamic interpretation relating directly to the Clausius
 302 definition of entropy. However, in this case we need to make assumptions stronger than (23) and
 303 (25). Specifically, we assume (24) and (26) which do not necessarily entail (23) and (25) respectively.
 304 For example, setting the conditional probability $p(\mathbf{x}_{n+1} | \mathbf{x}_n) = \frac{1}{Z_1} e^{(S-S_0)/k}$ does not mean that
 305 $W_1 = e^{(S-S_0)/k}$ is the transition probability.

Under the new, stronger, assumptions the conditional entropies can be related to the heat transferred
 in the transition, per temperature. Specifically, assumption (24) entails

$$h(\mathbf{x}_{n+1} | \mathbf{x}_n) = \log_2 Z_1 - \frac{1}{k \log 2} (S(\mathbf{x}_{n+1}) - S(\mathbf{x}_n)) = \log_2 Z_1 - \frac{1}{k \log 2} \int_{\mathbf{x}_n}^{\mathbf{x}_{n+1}} dq_{rev}/T. \quad (44)$$

306 where the last step used the definition of Clausius entropy (16). As per (16), this quantity is the same for
 307 all reversible paths between the past and next states. An example illustrating the transition ($\mathbf{x}_n \rightarrow \mathbf{x}_{n+1}$)
 308 can be given by a simple thermal system \mathbf{x}_n that is connected to a heat bath — that is, to a system in
 309 contact with a source of energy at temperature T . When the system X reaches a (new) equilibrium, e.g.,
 310 the state \mathbf{x}_{n+1} , due to its connection to the heat bath, the local conditional entropy $h(\mathbf{x}_{n+1} | \mathbf{x}_n)$ of the
 311 transition undergone by system X represents the heat transferred in the transition, per temperature.

Similarly, assumption (26) leads to

$$h(\mathbf{x}_{n+1} | \mathbf{x}_n, \mathbf{y}_n) = \log_2 Z_2 - \frac{1}{k \log 2} (\sigma_y) = \log_2 Z_2 - \frac{1}{k \log 2} \int_{\mathbf{x}_n \xrightarrow{\mathbf{y}_n} \mathbf{x}_{n+1}} dq/T, \quad (45)$$

312 where $\mathbf{x}_n \xrightarrow{\mathbf{y}_n} \mathbf{x}_{n+1}$ is the new path between \mathbf{x}_n and \mathbf{x}_{n+1} brought about by \mathbf{y}_n , and the entropy produced
 313 along this path is σ_y . That is, the first and the last points of the path over which we integrate heat transfers
 314 per temperature are unchanged but the path is affected by the source \mathbf{y} . This can be illustrated by a
 315 modified thermal system, still at temperature T but with heat flowing through some thermal resistance
 316 Y , while the system X repeats its transition from \mathbf{x}_n to \mathbf{x}_{n+1} .

Transfer entropy captures the difference between expressions (44) and (45), i.e., between the relevant
 amounts of heat transferred to the system X , per temperature.

$$t_{Y \rightarrow X}(n+1) = \log_2 \frac{Z_1}{Z_2} + \frac{1}{k \log 2} \left(\int_{\mathbf{x}_n \xrightarrow{\mathbf{y}_n} \mathbf{x}_{n+1}} dq/T - \int_{\mathbf{x}_n}^{\mathbf{x}_{n+1}} dq_{rev}/T \right). \quad (46)$$

317 Assuming that $Z_1 \approx Z_2$ is realistic, e.g. for quasistatic processes, then the additive constant disappears
 318 as well.

It is clear that if the new path is still reversible (e.g., when the thermal resistance is zero) then the
 source \mathbf{y} has not affected the resultant entropy change and we must have

$$\int_{\mathbf{x}_n}^{\mathbf{x}_{n+1}} dq_{rev}/T = \int_{\mathbf{x}_n \xrightarrow{\mathbf{y}_n} \mathbf{x}_{n+1}} dq/T \quad (47)$$

319 and $t_{Y \rightarrow X}(n+1) = 0$. This obviously occurs if and only if the source Y satisfies the independence
 320 condition (6), making the transfer entropy (46) equal to zero. In other words, we may again observe that if
 321 the local conditional entropy $h(\mathbf{x}_{n+1} | \mathbf{x}_n)$ corresponds to the resultant entropy change of the transition,
 322 then X and Y are dependent only when the external source Y , captured in the local conditional entropy
 323 $h(\mathbf{x}_{n+1} | \mathbf{x}_n, \mathbf{y}_n)$, brings about an irreversible internal change. If, however, the source Y changed the
 324 path in such a way that the process became irreversible, then $t_{Y \rightarrow X}(n+1) \neq 0$.

Finally, according to (19) and (20), the difference between the relevant heats transferred is $\int dq/T$,
 where q represents the heat flow received by the system from the exterior via the source Y , and hence

$$t_{Y \rightarrow X}(n+1) = \log_2 \frac{Z_1}{Z_2} - \frac{1}{k \log 2} \int dq/T. \quad (48)$$

325 In other words, local transfer entropy is proportional to the heat received or dissipated by the system
 326 from/to the exterior.

327 5. Causal effect: thermodynamic interpretation?

328 In this section we shall demonstrate that a similar treatment is not possible in general for causal
 329 effect. Again, we begin by considering local causal effect (15) of the source \mathbf{y}_n on destination \mathbf{x}_{n+1} ,
 330 while selecting s as the destination's past state \mathbf{x}_n :

$$f(\mathbf{y}_n \rightarrow \mathbf{x}_{n+1} | \hat{\mathbf{x}}_n) = \log_2 \frac{p(\mathbf{x}_{n+1} | \hat{\mathbf{y}}_n, \hat{\mathbf{x}}_n)}{\sum_{\mathbf{y}'_n} p(\mathbf{y}'_n | \hat{\mathbf{x}}_n) p(\mathbf{x}_{n+1} | \hat{\mathbf{y}}'_n, \hat{\mathbf{x}}_n)}. \quad (49)$$

331 Let us first consider conditions under which this representation reduces to the local transfer entropy.
 332 As pointed out by Lizier and Prokopenko [41], there are several conditions for such a reduction.

333 Firstly, \mathbf{y}_n and \mathbf{x}_n must be the only causal contributors to \mathbf{x}_{n+1} . In a thermodynamic setting, this
 334 means that there are no other sources affecting the transition from \mathbf{x}_n to \mathbf{x}_{n+1} , apart from \mathbf{y}_n .

Whenever this condition is met, and in addition, the combination $(\mathbf{y}_n, \mathbf{x}_n)$ is observed, it follows that

$$p(\mathbf{x}_{n+1} | \hat{\mathbf{y}}_n, \hat{\mathbf{x}}_n) = p(\mathbf{x}_{n+1} | \mathbf{y}_n, \mathbf{x}_n), \quad (50)$$

335 simplifying the numerator of Eq. (49).

Furthermore, there is another condition :

$$p(\mathbf{y}_n | \hat{\mathbf{x}}_n) \equiv p(\mathbf{y}_n | \mathbf{x}_n). \quad (51)$$

336 For example, it is met when the source \mathbf{y}_n is both causally and conditionally independent of the
 337 destination's past \mathbf{x}_n . Specifically, causal independence means $p(\mathbf{y}_n) \equiv p(\mathbf{y}_n | \hat{\mathbf{x}}_n)$, while conditional
 338 independence is simply $p(\mathbf{y}_n) \equiv p(\mathbf{y}_n | \mathbf{x}_n)$. Intuitively, the situation of causal and conditional
 339 independence means that inner workings of the system X under consideration do not interfere with
 340 the source Y . Alternatively, if X is the only causal influence on Y , the condition (51) also holds, as Y
 341 is perfectly "explained" by X , whether X is observed or imposed on. In general, though, the condition
 342 (51) means that the probability of \mathbf{y}_n if we impose a value $\hat{\mathbf{x}}_n$, is the same as if we had simply observed
 343 the value $\mathbf{x}_n = \hat{\mathbf{x}}_n$ without imposing in the system X .

Under the conditions (50) and (51), the denominator of Eq. (49) reduces to $p(\mathbf{x}_{n+1} | \mathbf{x}_n)$, yielding the equivalence between local causal effect and local transfer entropy

$$f(\mathbf{y}_n \rightarrow \mathbf{x}_{n+1} | \hat{\mathbf{x}}_n) = t_{Y \rightarrow X}(n+1). \quad (52)$$

344 In this case, the thermodynamic interpretation of transfer entropy would be applicable to causal effect as
345 well.

Whenever one of these conditions is not met, however, the reduction fails. Consider, for instance, the case when the condition (51) is satisfied, but the condition (50) is violated. For example, we may assume that there is some hidden source affecting the transition to \mathbf{x}_{n+1} . In this case, the denominator of Eq. (49) does not simplify much, and the component which may have corresponded to the entropy rate of the transition between \mathbf{x}_n and \mathbf{x}_{n+1} becomes

$$\log_2 \sum_{\mathbf{y}'_n} p(\mathbf{y}'_n | \mathbf{x}_n) p(\mathbf{x}_{n+1} | \hat{\mathbf{y}}'_n, \hat{\mathbf{x}}_n). \quad (53)$$

346 The interpretation of this irreducible component is important: presence of the imposed term $\hat{\mathbf{y}}'_n$ means
347 that one should estimate individual contributions of all possible states \mathbf{y} of the source Y , while varying
348 (i.e., imposing on) the state \mathbf{x}_n . This procedure becomes necessary because, in order to estimate the
349 causal effect of source \mathbf{y} , *in presence of some other hidden source*, one needs to check all possible
350 impositions on the source state \mathbf{y} . The terms of the sum under the logarithm in (53) inevitably vary in
351 their specific contribution, and so the sum cannot be analytically expressed as a single product under the
352 logarithm. This means that we cannot construct a direct thermodynamic interpretation of causal effect in
353 the same way that we did for the transfer entropy.

354 6. Discussion and Conclusions

355 In this paper we proposed a thermodynamic interpretation of transfer entropy: an information-
356 theoretic measure introduced by Schreiber [1] as the average information contained in the source about
357 the next state of the destination in the context of what was already contained in the destinations past. In
358 doing so we used a specialised Boltzmanns principle. This in turn produced a representation of transfer
359 entropy $t_{Y \rightarrow X}(n+1)$ as a difference of two entropy rates: one rate for a resultant transition within the
360 system of interest X and another rate for a possibly irreversible transition within the system affected
361 by an addition source Y . This difference was further shown to be proportional to the external entropy
362 production, Δ_{ext} , attributed to the source of irreversibility Y .

363 At this stage we would like to point out a difference between our main result, $t_{Y \rightarrow X}(n+1) \propto$
364 $-\Delta_{ext}$, and a representation for entropy production discussed by Parrondo et al. [22]. The latter work
365 characterised the entropy production in the total device, in terms of relative entropy, the Kullback-Leibler
366 divergence between the probability density ρ in phase space of some forward process and the probability
367 density $\tilde{\rho}$ of the corresponding and suitably defined time-reversed process. The consideration of Parrondo
368 et al. [22] does not involve any additional sources Y , and so transfer entropy is outside of the scope of
369 their study. Their main result characterised entropy production as $k d_{KL}(\rho || \tilde{\rho})$ which is equal to the
370 total entropy change in the total device. In contrast, in our study we consider the system of interest X
371 specifically, and characterise various entropy rates of X , but in doing so compare how these entropy

372 rates are affected by some source of irreversibility Y . In short, transfer entropy is shown to concur with
373 the entropy produced/dissipated by the system attributed to the external source Y .

374 We also briefly considered a case of fluctuations in the system X near an equilibrium, relating transfer
375 entropy to the difference in stabilities of the equilibrium, with respect to two scenarios: a default case and
376 the case with an additional source Y . This comparison was carried out with Kullback-Leibler divergences
377 of the corresponding transition probabilities.

378 Finally, we demonstrated that such a thermodynamic treatment is not applicable to information flow:
379 a measure introduced by Ay and Polani [18] in order to capture a causal effect. We argue that the
380 main reason is the interventional approach adopted in the definition of causal effect. We identified
381 several conditions ensuring certain dependencies between the involved variables, and showed that the
382 causal effect may also be interpreted thermodynamically — but in this case it reduces to transfer entropy
383 anyway. The highlighted difference once more shows a fundamental difference between transfer entropy
384 and causal effect: the former has a thermodynamic interpretation relating to the source of irreversibility
385 Y , while the latter is a construct that in general assumes an observer intervening in the system in a
386 particular way.

387 We hope that the proposed interpretation will further advance studies relating information theory and
388 thermodynamics, both in equilibrium and non-equilibrium settings, reversible and irreversible scenarios,
389 average and local scopes, etc.

390

391 Acknowledgements

392 The Authors are thankful to Ralf Der (Max Planck Institute for Mathematics in the Sciences, Leipzig)
393 who suggested and co-developed the example in subsection 3.6, and anonymous reviewers whose
394 suggestions significantly improved the paper.

395 References

- 396 1. Schreiber, T. Measuring Information Transfer. *Physical Review Letters* **2000**, *85*, 461–464.
- 397 2. Jakubowski, M.H.; Steiglitz, K.; Squier, R. Information transfer between solitary waves in the
398 saturable Schrödinger equation. *Physical Review E* **1997**, *56*, 7267.
- 399 3. Baek, S.K.; Jung, W.S.; Kwon, O.; Moon, H.T. Transfer Entropy Analysis of the Stock Market,
400 2005. arXiv:physics/0509014v2.
- 401 4. Moniz, L.J.; Cooch, E.G.; Ellner, S.P.; Nichols, J.D.; Nichols, J.M. Application of information
402 theory methods to food web reconstruction. *Ecological Modelling* **2007**, *208*, 145–158.
- 403 5. Chávez, M.; Martinerie, J.; Le Van Quyen, M. Statistical assessment of nonlinear causality:
404 application to epileptic EEG signals. *Journal of Neuroscience Methods* **2003**, *124*, 113–128.
- 405 6. Pahle, J.; Green, A.K.; Dixon, C.J.; Kummer, U. Information transfer in signaling pathways: a
406 study using coupled simulated and experimental data. *BMC Bioinformatics* **2008**, *9*, 139.
- 407 7. Lizier, J.T.; Prokopenko, M.; Zomaya, A.Y. Local information transfer as a spatiotemporal filter
408 for complex systems. *Physical Review E* **2008**, *77*, 026110.
- 409 8. Lizier, J.T.; Prokopenko, M.; Zomaya, A.Y. Information modification and particle collisions in
410 distributed computation. *Chaos* **2010**, *20*, 037109.

- 411 9. Lizier, J.T.; Prokopenko, M.; Zomaya, A.Y. Coherent information structure in complex
412 computation. *Theory in Biosciences* **2012**, *131*, 193–203.
- 413 10. Lizier, J.T.; Prokopenko, M.; Zomaya, A.Y. Local measures of information storage in complex
414 distributed computation. *Information Sciences* **2012**, *208*, 39–54.
- 415 11. Lizier, J.T.; Prokopenko, M.; Tanev, I.; Zomaya, A.Y. Emergence of Glider-like Structures
416 in a Modular Robotic System. Proceedings of the Eleventh International Conference on the
417 Simulation and Synthesis of Living Systems (ALife XI), Winchester, UK; Bullock, S.; Noble, J.;
418 Watson, R.; Bedau, M.A., Eds.; MIT Press: Cambridge, MA, 2008; pp. 366–373.
- 419 12. Lizier, J.T.; Prokopenko, M.; Zomaya, A.Y. The information dynamics of phase transitions
420 in random Boolean networks. Proceedings of the Eleventh International Conference on the
421 Simulation and Synthesis of Living Systems (ALife XI), Winchester, UK; Bullock, S.; Noble, J.;
422 Watson, R.; Bedau, M.A., Eds.; MIT Press: Cambridge, MA, 2008; pp. 374–381.
- 423 13. Lizier, J.T.; Pritam, S.; Prokopenko, M. Information dynamics in small-world Boolean networks.
424 *Artificial Life* **2011**, *17*, 293–314.
- 425 14. Lizier, J.T.; Heinzle, J.; Horstmann, A.; Haynes, J.D.; Prokopenko, M. Multivariate
426 information-theoretic measures reveal directed information structure and task relevant changes
427 in fMRI connectivity. *Journal of Computational Neuroscience* **2011**, *30*, 85–107.
- 428 15. Wang, X.R.; Miller, J.M.; Lizier, J.T.; Prokopenko, M.; Rossi, L.F. Quantifying and Tracing
429 Information Cascades in Swarms. *PLoS ONE* **2012**, *7*, e40084.
- 430 16. Lizier, J.T.; Prokopenko, M.; Cornforth, D.J. The information dynamics of cascading failures
431 in energy networks. Proceedings of the European Conference on Complex Systems (ECCS),
432 Warwick, UK, 2009, p. 54. ISBN: 978-0-9554123-1-8.
- 433 17. Kaiser, A.; Schreiber, T. Information transfer in continuous processes. *Physica D* **2002**, *166*, 43–
434 62.
- 435 18. Ay, N.; Polani, D. Information Flows in Causal Networks. *Advances in Complex Systems* **2008**,
436 *11*, 17–41.
- 437 19. Bennett, C.H. Notes on Landauer’s principle, reversible computation, and Maxwell’s Demon.
438 *Studies in History and Philosophy of Science Part B* **2003**, *34*, 501–510.
- 439 20. Piechocinska, B. Information erasure. *Physical Review A* **2000**, *61*, 062314.
- 440 21. Lloyd, S. *Programming the Universe*; Vintage Books: New York, 2006.
- 441 22. Parrondo, J.M.R.; den Broeck, C.V.; Kawai, R. Entropy production and the arrow of time. *New*
442 *Journal of Physics* **2009**, *11*, 073008.
- 443 23. Prokopenko, M.; Lizier, J.T.; Obst, O.; Wang, X.R. Relating Fisher information to order
444 parameters. *Physical Review E* **2011**, *84*, 041116.
- 445 24. Landauer, R. Irreversibility and heat generation in the computing process. *IBM Journal of*
446 *Research and Development* **1961**, *5*, 183–191.
- 447 25. Maroney, O.J.E. Generalizing Landauer’s principle. *Physical Review E* **2009**, *79*, 031105.
- 448 26. Jaynes, E.T. Information Theory and Statistical Mechanics. *Phys. Rev.* **1957**, *106*, 620–630.
- 449 27. Jaynes, E.T. Information Theory and Statistical Mechanics. II. *Phys. Rev.* **1957**, *108*, 171–190.
- 450 28. Crooks, G. Measuring Thermodynamic Length. *Physical Review Letters* **2007**, *99*, 100602+.

- 451 29. Liang, X.S. Information flow within stochastic dynamical systems. *Physical Review E* **2008**,
452 78, 031113.
- 453 30. Lüdtke, N.; Panzeri, S.; Brown, M.; Broomhead, D.S.; Knowles, J.; Montemurro, M.A.; Kell,
454 D.B. Information-theoretic sensitivity analysis: a general method for credit assignment in
455 complex networks. *Journal of The Royal Society Interface* **2008**, 5, 223–235.
- 456 31. Auletta, G.; Ellis, G.F.R.; Jaeger, L. Top-down causation by information control: from a
457 philosophical problem to a scientific research programme. *Journal of The Royal Society Interface*
458 **2008**, 5, 1159–1172.
- 459 32. Hlaváčková-Schindler, K.; Paluš, M.; Vejmelka, M.; Bhattacharya, J. Causality detection based
460 on information-theoretic approaches in time series analysis. *Physics Reports* **2007**, 441, 1–46.
- 461 33. Lungarella, M.; Ishiguro, K.; Kuniyoshi, Y.; Otsu, N. Methods for quantifying the causal
462 structure of bivariate time series. *International Journal of Bifurcation and Chaos* **2007**,
463 17, 903–921.
- 464 34. Ishiguro, K.; Otsu, N.; Lungarella, M.; Kuniyoshi, Y. Detecting direction of causal interactions
465 between dynamically coupled signals. *Physical Review E* **2008**, 77, 026216.
- 466 35. Sumioka, H.; Yoshikawa, Y.; Asada, M. Causality Detected by Transfer Entropy Leads
467 Acquisition of Joint Attention. Proceedings of the 6th IEEE International Conference on
468 Development and Learning (ICDL 2007), London. IEEE, 2007, pp. 264–269.
- 469 36. Vejmelka, M.; Palus, M. Inferring the directionality of coupling with conditional mutual
470 information. *Physical Review E* **2008**, 77, 026214.
- 471 37. Verdes, P.F. Assessing causality from multivariate time series. *Physical Review E* **2005**,
472 72, 026222–9.
- 473 38. Tung, T.Q.; Ryu, T.; Lee, K.H.; Lee, D. Inferring Gene Regulatory Networks from Microarray
474 Time Series Data Using Transfer Entropy. Proceedings of the Twentieth IEEE International
475 Symposium on Computer-Based Medical Systems (CBMS '07), Maribor, Slovenia; Kokol, P.;
476 Podgorelec, V.; Mičetič-Turk, D.; Zorman, M.; Verlič, M., Eds.; IEEE: Los Alamitos, USA,
477 2007; pp. 383–388.
- 478 39. Van Dijck, G.; Van Vaerenbergh, J.; Van Hulle, M.M. Information Theoretic Derivations for
479 Causality Detection: Application to Human Gait. Proceedings of the International Conference
480 on Artificial Neural Networks (ICANN 2007), Porto, Portugal; Sá, J.M.d.; Alexandre, L.A.;
481 Duch, W.; Mandic, D., Eds.; Springer-Verlag: Berlin/Heidelberg, 2007; Vol. 4669, *Lecture Notes*
482 *in Computer Science*, pp. 159–168.
- 483 40. Hung, Y.C.; Hu, C.K. Chaotic Communication via Temporal Transfer Entropy. *Physical Review*
484 *Letters* **2008**, 101, 244102.
- 485 41. Lizier, J.T.; Prokopenko, M. Differentiating information transfer and causal effect. *European*
486 *Physical Journal B* **2010**, 73, 605–615.
- 487 42. Wuensche, A. Classifying cellular automata automatically: Finding gliders, filtering, and relating
488 space-time patterns, attractor basins, and the Z parameter. *Complexity* **1999**, 4, 47–66.
- 489 43. Solé, R.V.; Valverde, S. Information transfer and phase transitions in a model of internet traffic.
490 *Physica A* **2001**, 289, 595–605.

- 491 44. Solé, R.V.; Valverde, S. Information Theory of Complex Networks: On Evolution and
492 Architectural Constraints. In *Complex Networks*; Ben-Naim, E.; Frauenfelder, H.; Toroczkai,
493 Z., Eds.; Springer: Berlin / Heidelberg, 2004; Vol. 650, *Lecture Notes in Physics*, pp. 189–207.
- 494 45. Prokopenko, M.; Boschiatti, F.; Ryan, A.J. An Information-Theoretic Primer on Complexity,
495 Self-Organization, and Emergence. *Complexity* **2009**, *15*, 11–28.
- 496 46. Takens, F. Detecting strange attractors in turbulence. In *Dynamical Systems and Turbulence*,
497 *Warwick 1980*; Rand, D.; Young, L.S., Eds.; Lecture Notes in Mathematics, Springer: Berlin /
498 Heidelberg, 1981; pp. 366–381.
- 499 47. Granger, C.W.J. Investigating causal relations by econometric models and cross-spectral
500 methods. *Econometrica* **1969**, *37*, 424–438.
- 501 48. Fano, R. *Transmission of Information: A Statistical Theory of Communications*; The MIT Press:
502 Cambridge, MA, 1961.
- 503 49. Manning, C.D.; Schütze, H. *Foundations of Statistical Natural Language Processing*; The MIT
504 Press: Cambridge, MA, USA, 1999.
- 505 50. Dasan, J.; Ramamohan, T.R.; Singh, A.; Nott, P.R. Stress fluctuations in sheared Stokesian
506 suspensions. *Physical Review E* **2002**, *66*, 021409.
- 507 51. MacKay, D.J. *Information Theory, Inference, and Learning Algorithms*; Cambridge University
508 Press: Cambridge, 2003.
- 509 52. Pearl, J. *Causality: Models, Reasoning, and Inference*; Cambridge University Press: Cambridge,
510 2000.
- 511 53. Goyal, P. Information Physics – Towards a New Conception of Physical Reality. *Information*
512 **2012**, *3*, 567–594.
- 513 54. Sethna, J.P. *Statistical mechanics: entropy, order parameters, and complexity*; Oxford University
514 Press: Great Clarendon Street, Oxford OX2 6DP, 2006.
- 515 55. Seife, C. *Decoding the universe*; Penguin Group: New York, 2006.
- 516 56. Einstein, A. Über einen die Erzeugung und Verwandlung des Lichtes betreffenden heuristischen
517 Gesichtspunkt. *Annals of Physics* **1905**, *322*, 132–148.
- 518 57. Norton, J.D. Atoms, Entropy, Quanta: Einstein’s Miraculous Argument of 1905. *Studies in*
519 *History and Philosophy of Modern Physics* **2006**, *37*, 71–100.
- 520 58. Barnett, L.; Buckley, C.L.; Bullock, S. Neural complexity and structural connectivity. *Phys. Rev.*
521 *E* **2009**, *79*, 051914.
- 522 59. Ay, N.; Bernigau, H.; Der, R.; Prokopenko, M. Information-driven self-organization: the
523 dynamical system approach to autonomous robot behavior. *Theory in Biosciences* **2012**,
524 *131*, 161–179.
- 525 60. Evans, D.J.; Cohen, E.G.D.; Morriss, G.P. Probability of second law violations in shearing steady
526 states. *Physical Review Letters* **1993**, *71*, 2401–2404.
- 527 61. Searles, D.J.; Evans, D.J. Fluctuations Relations for Nonequilibrium Systems. *Australian*
528 *Journal of Chemistry* **2004**, *57*, 1129–1123.
- 529 62. Crooks, G.E. Entropy production fluctuation theorem and the nonequilibrium work relation for
530 free energy differences. *Physical review E* **1999**, *60*, 2721–2726.

- 531 63. Jarzynski, C. Nonequilibrium work relations: foundations and applications. *The European*
532 *Physical Journal B - Condensed Matter and Complex Systems* **2008**, *64*, 331–340.
- 533 64. Schlögl, F. Information measures and thermodynamic criteria for motion. In *Structural stability*
534 *in physics: Proceedings of two International Symposia on Applications of Catastrophe Theory*
535 *and Topological Concepts in Physics, Tübingen, May 2-6 and December 11-14, 1978*; Güttinger,
536 W.; Eikemeier, H., Eds.; Springer, 1979; Vol. 4, *Springer series in synergetics*, pp. 199–209.

537 © February 6, 2013 by the authors; submitted to *Entropy* for possible open access
538 publication under the terms and conditions of the Creative Commons Attribution license
539 <http://creativecommons.org/licenses/by/3.0/>.