

The problem of infinite information flow*

Zheng Bian[†] and Erik M. Bollt[‡]

Abstract. We study conditional mutual information (cMI) between a pair of variables X, Y given a third one Z and derived quantities including transfer entropy (TE) and causation entropy (CE) in the dynamically relevant context where $X = T(Y, Z)$ is determined by Y, Z via a deterministic transformation T . Under mild continuity assumptions on their distributions, we prove a zero-infinity dichotomy for cMI for a wide class of T , which gives a yes-or-no answer to the question of information flow as quantified by TE or CE. Such an answer fails to distinguish between the relative amounts of information flow. To resolve this problem, we propose a discretization strategy and a conjectured formula to discern the *relative ambiguities* of the system, which can serve as a reliable proxy for the relative amounts of information flow. We illustrate and validate this approach with numerical evidence.

Key words. information flow, causal inference, information theory, entropy, Kullback–Leibler divergence, mutual information, conditional mutual information

MSC codes. 94A17

1. Introduction. Quantifying information flow is a critical task for understanding complex systems in various scientific disciplines, from neuroscience [26, 25, 20] to financial markets [8, 3]. Information measures such as mutual information (MI), conditional mutual information (cMI) [7], transfer entropy (TE) [19], and causation entropy (CE) [23], have become essential tools for this purpose.

Tracing back to the classic Weiner-Granger causality [11, 12, 4, 15], a central idea that underlies these information theoretic methods of quantifying information flow is the notion of *disambiguation* in a predictive framework. In contrast to the experimentalist approach, which infers causality from outcomes of perturbations and experiments, the predictive framework, which we consider below, is premised on alternative formulations of the forecasting question, with and without considering the influence of an external system.

Formulated by Schreiber [19] in 2000, TE is a quantitative attempt in this predictive framework. We think of $V = \{V_t\}$ and $U = \{U_t\}$ as stochastic processes indexed by discrete time $t = 0, 1, \dots$; for a concrete example, imagine that V, U record EEG times series data from different parts of the brain. We expect that the present state V_t informs about the future state V_{t+1} and are interested in determining whether the present state U_t also informs about V_{t+1} . If V_{t+1} is conditionally independent of U_t given V_t , then the knowledge about the state of U_t does not resolve any uncertainty about the state of V_{t+1} , assuming one already has access to the state of V_t . In this case, we would like to conclude no information flow from U to V at

*Submitted to the editors March 2025.

Funding: E.M.B. and Z.B. are supported by the NSF-NIH-CRCNS, and E.M.B. is also supported by DARPA RSDN, the ARO, and the ONR.

[†]Clarkson Center for Complex Systems Science (C³S²), Potsdam, NY 13699 USA, (zheng@bian-zheng.cn).

[‡]Department of Electrical and Computer Engineering, Clarkson University, Potsdam, NY 13699 USA, and Clarkson Center for Complex Systems Science (C³S²), Potsdam, NY 13699 USA, (ebollt@clarkson.edu).

36 time t and zero TE accordingly. Otherwise, any deviation from this conditional independence
 37 indicates the presence of information flow, to be captured and quantified by some positive
 38 value of TE measured in bits per time unit.

39 By a slight generalization of Schreiber’s original formulation and in agreement with the
 40 usual definition for discrete variables, we define TE

$$41 \quad (1.1) \quad T_{U \rightarrow V, t} := I(V_{t+1}; U_t | V_t)$$

42 to be the conditional mutual information of V_{t+1}, U_t given V_t . For simplicity, this is the case of
 43 lag length 1; longer lags are allowed in general. Causation entropy, proposed by Sun, Taylor
 44 and Bollt [23], generalizes TE to infer network connectivity [21, 22, 1, 16], by also building
 45 in conditioning on ternary influences as a way to resolve the differences between direct and
 46 indirect interactions. The precise definition of cMI will be given in Section 2. Roughly speak-
 47 ing, it quantifies the deviation from conditional independence of a pair of random variables
 48 conditioned on a third variable.

49 **1.1. Zero-infinity dichotomy.** Consider a typical situation from dynamical systems, where
 50 the random variable V_{t+1} is determined by U_t, V_t via some deterministic map T , that is,

$$51 \quad (1.2) \quad V_{t+1} = T(U_t, V_t).$$

52 If V_{t+1} does not depend on U_t , that is, $V_{t+1} = T_0(V_t)$, then we trivially have zero information
 53 flow $T_{U \rightarrow V, t} = 0$. In terms of probability distributions, this case corresponds to the regular
 54 conditional probability $\mathbb{P}(V_{t+1} \in \cdot | V_t = v_t) = \delta_{T_0(v_t)}$ being a dirac delta.

55 Otherwise, one expects $T_{U \rightarrow V, t} > 0$ to quantify the amount of information flowing from U
 56 to V at time t . For example, if the map T is highly “ambiguous”, then the knowledge about
 57 the states of U_t, V_t does not resolve much uncertainty about the state of V_{t+1} .

58 *Example 1.1.* Consider two maps $T_1(u, v) = 100(u + v) \bmod 1$ and $T_2(u, v) = u + v$
 59 $\bmod 1$. The knowledge about the states of U_t, V_t up to 10^{-2} precision is completely lost via
 60 T_1 and trivially informs that $V_{t+1} = T_1(U_t, V_t)$ lies in $[0, 1]$, whereas this knowledge under T_2
 61 informs about the state of $V_{t+1} = T_2(U_t, V_t)$ up to precision 2×10^{-2} . Therefore, we may
 62 expect $T_{U \rightarrow V, t}$ to be smaller in the more ambiguous case of $V_{t+1} = T_1(U_t, V_t)$ than in the case
 63 of $V_{t+1} = T_2(U_t, V_t)$.

64 However, under some mild continuity assumptions on the distribution of V_{t+1} , we see that in
 65 both cases, $T_{U \rightarrow V, t} = \infty$. This holds more generally for any measurable map T . Throughout
 66 this paper, we assume that the random variables take values in standard measurable spa-
 67 ces, unless otherwise stated. This implies the existence and essential uniqueness of regular
 68 conditional probabilities and disintegrations; for details see Appendix A.

69 **Theorem A (infinite information flow):** *Assume that for a positive measure set of*
 70 *outcomes v_t of V_t , the regular conditional probability distribution $\mathbb{P}(V_{t+1} \in \cdot | V_t = v_t)$ of V_{t+1}*
 71 *in Eq. (1.2) charges an atomless continuum. Then, the transfer entropy $T_{U \rightarrow V, t}$ from U to V*
 72 *at time t is infinite.*

73 *Remark 1.2.* The positive measure set is with respect to the distribution of V_t . We say
 74 that a probability measure μ charges an atomless continuum if there is a measurable set B

75 such that $\mu(B) > 0$ and $\mu(\{b\}) = 0$ for each point $b \in B$. The assumption of Theorem A
 76 says that V_t alone does not fully determine V_{t+1} but rather leaves a rich continuum of possible
 77 values for V_{t+1} . This is the case, for example, when $V_{t+1} = T_1(U_t, V_t)$ or $V_{t+1} = T_2(U_t, V_t)$ as
 78 in Example 1.1 with V_t and U_t independent and following the uniform distribution on $[0, 1]$.

79 Theorem 3.7 gives an equivalent but slightly different formulation of Theorem A and is
 80 proven in Section 3.3. The zero-infinity dichotomy of $T_{U \rightarrow V, t}$ gives a yes-or-no answer to the
 81 question of information flow.

82 A key step in the proof of Theorem A is to disintegrate the conditional mutual informa-
 83 tion into mutual information between conditioned variables. We believe that this result is
 84 interesting in its own right and state it below.

85 **Theorem B (disintegration of conditional mutual information):** *The conditional*
 86 *mutual information $I(X; Y|Z)$ of three random variables X, Y, Z is the average of the mutual*
 87 *information $I(X_z; Y_z)$ between conditioned versions X_z, Y_z of X, Y defined in Eq. (2.3), that*
 88 *is,*

$$89 \quad I(X; Y|Z) = \int I(X_z; Y_z) dP_Z(z).$$

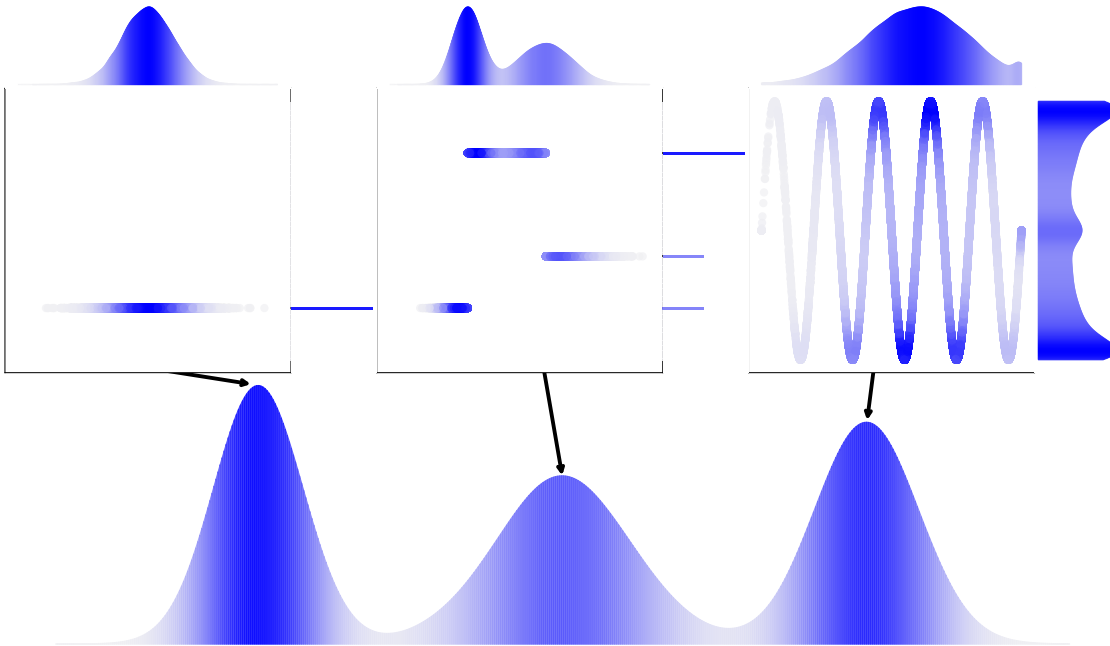


Figure 1. *Disintegrated distributions. The main histogram at the bottom illustrates the distribution P_Z of variable Z , which, together with Y , determines $X = T(Y, Z)$ via a measurable map T . The joint distribution P_{XYZ} disintegrates into $(P_{XYZ})_z$ for each realization of $Z = z$, which can be interpreted as the joint distribution $P_{X_z Y_z}$ of the conditioned versions X_z, Y_z of X, Y . The left, center and right subplots above the main histogram illustrate three typical disintegrated distributions $(P_{XYZ})_z = P_{X_z Y_z}$, where X_z follows a constant, atomic and continuous distribution, respectively. In each subplot, the scatter plot shows the joint distribution $P_{X_z Y_z}$, the top histogram shows the marginal distribution P_{Y_z} , and the right histogram shows the marginal distribution P_{X_z} . The intensity of the blue gradient indicates regions of high probability density.*

90 *Remark 1.3.* The conditioned variables X_z, Y_z describe the probabilistic landscape once
 91 the uncertainty about Z is removed, by assuming that the outcome of Z is z . This allows
 92 the intermediate measurement of $I(X_z; Y_z)$ on this particular outcome. By averaging across
 93 all outcomes of Z , the full conditional mutual information $I(X; Y|Z)$ is recovered. We illus-
 94 trate pictorially three typical scenarios in Figure 1; the subplots show the joint $P_{X_z Y_z}$ and
 95 marginal distributions P_{X_z}, P_{Y_z} of pairs of random variables X_z, Y_z above the main histogram
 96 illustrating the distribution of Z . Proposition 2.8 gives an equivalent but slightly different
 97 formulation of Theorem B and is proven in Section 2.3. The main technical step involves the
 98 proper construction of X_z, Y_z in Eq. (2.3) and the equivalence of disintegration and regular
 99 conditional probability in our context.

100 Theorem B reduces the analysis of TE or cMI in Theorem A to that of MI between con-
 101 ditioned variables. The exhaustive analysis of MI in the deterministic context thus completes
 102 the proof of Theorem A.

103 In practice, one computes TE from a finite amount of data and obtains finite positive
 104 values of $T_{U \rightarrow V, t}$. As noted in [6], much of the literature that applies TE to detect information
 105 flow focuses on establishing that $T_{U \rightarrow V, t}$ is statistically significantly different from zero, and
 106 treats the finite positive values of $T_{U \rightarrow V, t}$ as mere artifacts of finite sampling.

107 As discussed in Example 1.1, a more ambiguous map such as T_1 allows through less
 108 information flow, which should be reflected by a smaller value of $T_{U \rightarrow V, t}$. Of course, this
 109 intuitive assumption is valid for discrete variables. However, it lacks theoretical justification in
 110 the case of continuous variables as pointed out by Theorem A, which is typical for applications
 111 to dynamical systems. We refer to this discrepancy between the practically obtained finite TE
 112 values and the theoretic zero-infinity dichotomy as the *problem of infinite information flow*.

113 **1.2. Resolution by discretization.** In light of Theorem B, it suffices to analyze the pair-
 114 wise $I(X; Y)$ for $X = T(Y)$, seeing that $I(X; Y|Z)$ can be obtained by averaging across
 115 $I(X_z, Y_z)$ for pairs of conditioned variables X_z, Y_z . A resolution of the problem of infinite
 116 information flow needs to achieve two things:

- 117 (R1) modify the model so as to obtain finite values for $I(X; Y)$,
- 118 (R2) by comparing the relative values, distinguish between the relative amounts of infor-
 119 mation flow.

120 By adding white noise to the map T as employed in [24], one can easily achieve (R1) as a
 121 blurring effect. However, we will show in Appendix B that this strategy still falls short of (R2).
 122 In fact, we prove for Bernoulli maps with uniformly distributed additive noise of amplitude
 123 ϵ , uniformly distributed Y and hence X , the resulting finite value of $I(X; Y)$ is $\ln \frac{1}{\epsilon}$, which is
 124 a function of the noise amplitude alone, independent of the expanding rate of the Bernoulli
 125 map. In this sense, the addition of white noise does not achieve (R2) because the resulting
 126 finite values of $I(X; Y)$ cannot distinguish between the relative dynamical ambiguities of the
 127 Bernoulli systems.

128 We propose discretization as a strategy to achieve both (R1) and (R2) and illustrate in
 129 the one-dimensional case.

130 **Conjecture C (relative ambiguity of (T, Y)):** *Suppose that X, Y are \mathbb{R} -valued random*

131 variables with continuous probability density functions f_X, f_Y , respectively, and that there is a
 132 piecewise C^1 map T for which $|T'| \geq 1$ and $X = T(Y)$. Consider the discretization by uniform
 133 mesh of size $\Delta > 0$, that is,

$$134 \quad \Pi^\Delta : \mathbb{R} \rightarrow \mathbb{Z}\Delta, \quad (\Pi^\Delta)^{-1}\{i\Delta\} = [i\Delta, (i+1)\Delta), \quad i \in \mathbb{Z}.$$

135 Then, in the limit as $\Delta \rightarrow 0^+$, the discretized variables $X^\Delta := \Pi^\Delta X, Y^\Delta := \Pi^\Delta Y$ satisfy

$$136 \quad I(X^\Delta; Y^\Delta) + \ln \Delta \rightarrow H(X) - \int \ln |T'| f_Y dy =: -A_T(Y),$$

137 where $H(X) := -\int f_X \ln f_X dx$ is the differential entropy of X and the quantity $A_T(Y)$ shall
 138 be called the relative ambiguity of system (T, Y) .

139 *Remark 1.4.* In the special case of $T = \text{id}$, we have $I(X^\Delta; X^\Delta) + \ln \Delta \rightarrow H(X)$ and recover
 140 the relation between Shannon entropy and differential entropy, see e.g. [7, Section 9.3]. More
 141 generally, it is clear that in the refinement limit of the discretization, i.e., as $\Delta \rightarrow 0^+$, the
 142 MI between the discretized variables $I(X^\Delta, Y^\Delta)$ tends to the infinite theoretic value $I(X; Y)$.
 143 This is not our primary concern, however. What is more interesting is the behavior for finite
 144 Δ^{-1} . Namely, for any finite Δ^{-1} , the intuition that a more ambiguous system (T, Y) with
 145 large relative ambiguity $A_T(Y)$ allows through less information is reflected by a smaller value
 146 of $I(X^\Delta, Y^\Delta)$. In this sense, discretization achieves both (R1) and (R2), resolving the problem
 147 of infinite information flow.

148 Note that the relative ambiguity $A_T(Y)$ involves an entropy and an exponent, which nat-
 149 urally suggests a link to the Pesin entropy formula [18]. However, we defer further discussions
 150 on this link, as well as the proof and generalization of Conjecture C, to a separate ongoing
 151 work.

152 Below, we validate Conjecture C with numerical evidence in some concrete dynamical
 153 examples. A sketch of the derivation of the conjectured formula for $A_T(Y)$ is included in the
 154 Appendix C.

155 *Example 1.5 (Bernoulli interval maps).* Let the random variable $X = E_d(Y)$ be determined
 156 by Y via the piecewise linear expanding map $E_d : [0, 1] \rightarrow [0, 1]$, $d \in \mathbb{Z}$, $d \geq 2$, on the unit
 157 interval given by

$$158 \quad E_d(x) = d \cdot x \pmod{1}.$$

159 Assume Y follows a continuous distribution (we consider uniform and Gaussian $\mathcal{N}_{[0,1]}(0.3, 0.02)$
 160 centered at 0.3 with variance 0.02 truncated between 0 and 1) on the interval. By Theorem
 161 A, or more directly, Theorem 3.6, $I(X; Y) = \infty$.

162 From Conjecture C, we have zero differential entropy of the uniformly distributed variable
 163 X and a constant expansion rate $|T'| = d$, which yields $A_T(Y) = \ln d$.

164 A direct calculation, see Section 4.1, shows that if Y is uniformly distributed in $[0, 1]$, then
 165 so is X and

$$166 \quad I(X^\Delta; Y^\Delta) = \ln \Delta^{-1} - \ln d = -\ln \Delta + A_T(Y),$$

167 in agreement with Conjecture C.

168 In Figure 2, we set $\Delta^{-1} = 300$. The left and center panels show the scatter plots of
 169 the joint distribution $P_{X^\Delta Y^\Delta}$ of the discretized variables X^Δ, Y^Δ , together with the marginal
 170 distribution P_{Y^Δ} on the top and P_{X^Δ} on the right of the scatter plots. We take Y to follow
 171 the uniform distribution in the left panel in blue and the Gaussian $\mathcal{N}_{[0,1]}(0.3, 0.02)$ in the
 172 center panel in red. The intensity of the colors indicates the high probability density. The
 173 right panel shows the mutual information $I(X^\Delta, Y^\Delta)$ decreases as the expansion rate d of the
 174 Bernoulli map E_d increases. The blue and red dots correspond to the cases of Y following
 175 the E_d -invariant uniform distribution and the Gaussian $\mathcal{N}_{[0,1]}(0.3, 0.02)$, respectively. For
 176 comparison, we superimpose the Conjecture prediction $\ln \Delta^{-1} + H(X) - \ln d$ in dashed lines.
 177 Observe that the dots from empirical calculations fit well with the Conjecture C predic-
 178 tions in dashed lines in both the uniform and Gaussian cases. In comparison to the uniform
 179 distribution, the tight Gaussian distribution $\mathcal{N}_{[0,1]}(0.3, 0.02)$ of Y results in a smaller (in fact,
 180 negative) differential entropy term $H(X)$ and hence a bigger relative ambiguity $A_T(Y)$ of the
 181 system (T, Y) and a smaller discretized mutual information. As the Bernoulli expanding rate
 182 d increases, the system (T, Y) becomes more ambiguous in both the uniform and Gaussian
 183 cases, and hence $I(X^\Delta, Y^\Delta)$ decreases. For very large d , the expansion is so strong that even
 184 the tight Gaussian distribution of Y smoothens to an almost uniform distribution of X via
 185 E_d and we see convergence of the two curves. This example validates both Conjecture C and
 186 the discretization strategy's ability to achieve (R1-2).

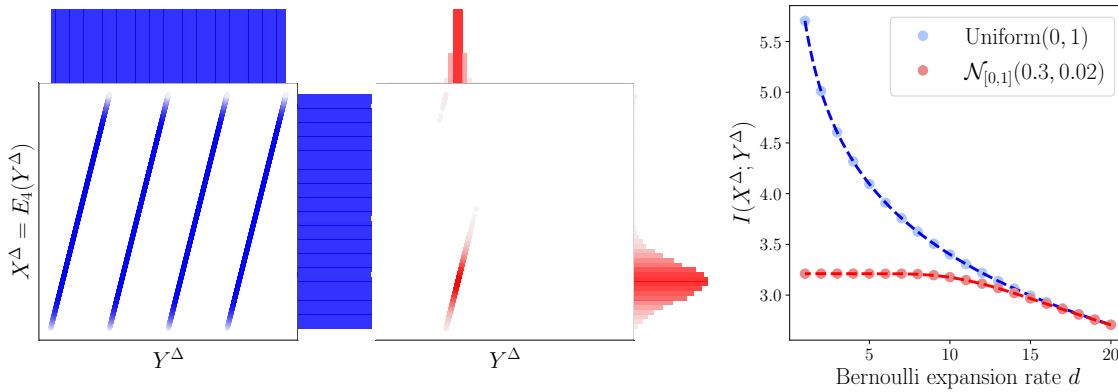


Figure 2. Discretization via uniform $\Delta^{-1} = 300$ partition of continuous random variable $X = E_d(Y)$ determined by variable Y via the Bernoulli map $E_d : x \mapsto d \cdot x \pmod{1}$. In the left and middle panels, the scatter plots show the joint distribution $P_{X^\Delta Y^\Delta}$ of the discretized variables X^Δ, Y^Δ , together with the marginal distributions P_{Y^Δ} at the top and P_{X^Δ} on the right. The blue and red plots correspond to Y following the uniform and Gaussian $\mathcal{N}_{[0,1]}(0.3, 0.02)$ distributions, respectively. Here, $\mathcal{N}_{[0,1]}(0.3, 0.02)$ means the Gaussian distribution centered at 0.3 with variance 0.02 and truncated between 0 and 1. The right panel plots for each Bernoulli expansion rate d , the corresponding $I(X^\Delta; Y^\Delta)$ of the discretized variables. The blue and red dots correspond to the empirical calculations of uniform and Gaussian $\mathcal{N}_{[0,1]}(0.3, 0.02)$ distributions, respectively. The dashed lines show the theoretic predictions from Conjecture C.

187 The next example illustrates the discretization strategy in a nonlinear case and beyond
 188 the scope of Conjecture C (because the map has contracting regions).

189 **Example 1.6 (Sine box functions).** Let the random variable $X = S_n(Y)$ be determined by

190 Y via the sine box function $S_n : [0, 1] \rightarrow [0, 1]$ given by

$$191 \quad S_n(x) := \frac{1 + \sin 2\pi nx}{2}, \quad n = 1, 2, \dots$$

192 We consider two continuous distributions for Y , namely, the uniform distribution and the
 193 absolutely continuous S_n -invariant probability (acip) distribution. The acip is approximated
 194 by a long trajectory $\{y_t\}$, $y_{t+1} = S_n(y_t)$, $t = \tau_0, \tau_0 + 1, \dots, \tau_0 + \tau - 1$ of length $\tau = 10^6$ with
 195 the first $\tau_0 = 1000$ iterates discarded as transients. In both cases, we have $I(X; Y) = \infty$ by
 196 Theorem A, or more directly, Theorem 3.6.

197 In Figure 3, we discretize X, Y the same way as in Example 1.5. For S_4 , we show the scatter
 198 plots of $P_{X^\Delta Y^\Delta}$ and histogram of P_{Y^Δ} at the top and P_{X^Δ} on the right of the left and center
 199 panels. The uniform P_Y shown in blue on the left is not invariant for S_n , but the red acip in the
 200 middle is S_n -invariant. The right panel shows that with Y following either uniform or acip
 201 distribution, the mutual information $I(X^\Delta; Y^\Delta)$ between the discretized variables X^Δ, Y^Δ
 202 decreases as the function S_n becomes more ambiguous (as n increases). The calculation and
 203 simulation details are presented in Section 4.2.

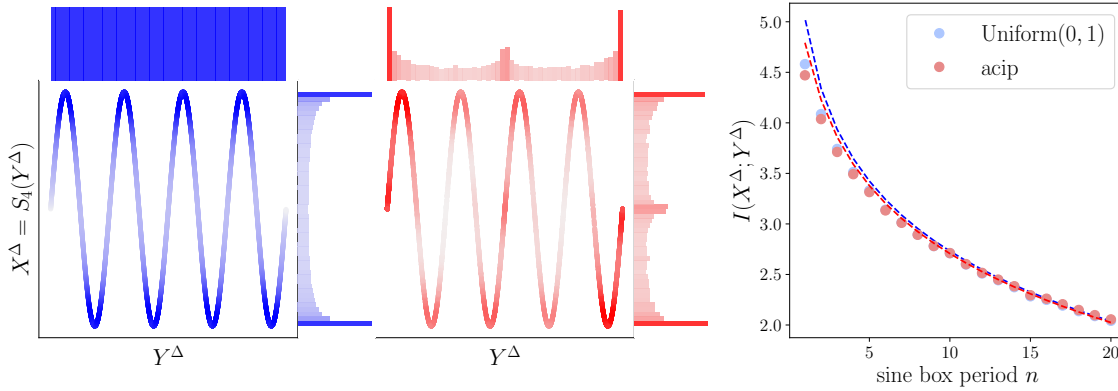


Figure 3. Discretization via uniform $\Delta^{-1} = 300$ partition of continuous random variable $X = S_n(Y)$ determined by variable Y via the sine box function $S_n : x \mapsto \frac{1 + \sin 2\pi nx}{2}$. In the left and middle panels, the scatter plots show the joint distribution $P_{X^\Delta Y^\Delta}$ of the discretized variables X^Δ, Y^Δ , together with the marginal distributions P_{Y^Δ} at the top and P_{X^Δ} on the right. The right panel plots for each n , the corresponding MI $I(X^\Delta; Y^\Delta)$ of the discretized variables, with the empirical values shown in dots and Conjectured values in dashed lines. The blue and red colors correspond to Y following the uniform and acip distributions, respectively.

204 We remark that the sine box example falls outside the scope of Conjecture C because
 205 S_n has contracting regions near $\frac{k}{2n} + \frac{1}{4n}$ for each $n = 1, 2, \dots$ and $k = 0, \dots, 2n - 1$, where
 206 our Conjectured formula fails. It turns out that these contracting regions are assigned a
 207 higher weight for smaller values of n and uniform and acip densities of f_Y , leading to a bigger
 208 discrepancy between the empirical and Conjectured $I(X^\Delta, Y^\Delta)$ values for small n . In spite of
 209 this, it is remarkable that our formula still captures the trend that as n increases, the relative
 210 ambiguity of S_n increases and $I(X^\Delta, Y^\Delta)$ decreases. This example illustrates the validity of
 211 the discretization strategy.

212 To obtain meaningful finite values of TE or cMI in Eq. (1.1) that can distinguish the
 213 relative amounts of information flow, we discretize each conditioned version Y_z and $X_z =$

214 $T_z(Y_z) = T(Y_z, z)$ to obtain meaningful finite values of MI $I(X_z; Y_z)$ as in Examples 1.5, 1.6,
 215 and then average/integrate across all versions of z against the marginal distribution P_Z (or
 216 its discretization) in the sense of Theorem B. The numerical computations of TE in practice,
 217 in our view, essentially implement a similar discretization scheme.

218 **Organization of the paper.** In Section 2 we review the definition and properties of MI
 219 and cMI and end with Proposition 2.8 to decompose cMI into disintegrated MI of conditioned
 220 versions of the original variables. In Section 3, we analyze the dichotomy properties of MI
 221 and cMI leading to the proof of the Theorem of infinite information flow. In Section 4,
 222 we present detailed calculations and simulations for the illustrative Bernoulli and sine box
 223 examples 1.5, 1.6. In the Appendix, we discuss the key technical results on standard spaces,
 224 regular conditional probability, disintegration, and the effect of additive white noise.

225 **Acknowledgments.** We thank Tiago Pereira and Edmilson Roque dos Santos for helpful
 226 discussions and comments. Z.B. and E.M.B. are supported by the NSF-NIH-CRCNS. E.M.B.
 227 is also supported by DARPA RSDN, the ARO, and the ONR.

228 **2. Background on cMI.** We review notions and properties of Kullback-Leibler divergence
 229 in Section 2.1, entropy and mutual information in Section 2.2, and the conditional mutual in-
 230 formation in Section 2.3. Some technical definitions and constructions, including the standard
 231 measurable space and regular conditional probability, are essential for the general definition
 232 of the conditional mutual information and therefore are also briefly reviewed in the Appendix.
 233 More details can be found in [13, 14].

234 Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and $f : \Omega \rightarrow A$ a measurable function (also called
 235 random variable) taking values in the measurable space (A, \mathcal{B}_A) called the *alphabet*. Denote
 236 the *distribution* of f on (A, \mathcal{B}_A) by

$$237 \quad P_f := f_*\mathbb{P}.$$

238 When A is a finite/countable set, we say that the alphabet is finite/discrete. For several
 239 random variables f_1, \dots, f_n , we denote their joint distribution by $P_{f_1 \dots f_n} = (f_1, \dots, f_n)_*\mathbb{P}$
 240 and the product measure of their marginal distributions by $P_{f_1} \otimes \dots \otimes P_{f_n} = ((f_1)_*\mathbb{P}) \otimes \dots \otimes$
 241 $((P_{f_n})_*\mathbb{P})$.

242 **2.1. Kullback-Leibler divergence.** First consider the special case where A is a finite set
 243 and $\mathcal{B}_A = 2^A$. Given two probability measures P, M on (A, \mathcal{B}_A) , the *Kullback-Leibler diver-*
 244 *gence of P with respect to M* is defined to be

$$245 \quad \text{KL}(P\|M) := \sum_{a \in A} P(a) \ln \frac{P(a)}{M(a)}.$$

246 Note that this makes sense only when $M(a) = 0$ implies $P(a) = 0$, i.e., $P \ll M$. In this case
 247 we define $0 \ln \frac{0}{0} := 0$; otherwise, $\text{KL}(P\|M)$ is defined to be ∞ .

248 Now consider the general case: two probability measures P, M on an arbitrary measurable
 249 space (Ω, \mathcal{F}) . The *Kullback-Leibler divergence* $\text{KL}(P\|M)$ of P with respect to M is defined
 250 as

$$251 \quad \text{KL}(P\|M) := \sup_f \text{KL}(P_f\|M_f),$$

252 where the supremum is taken over all random variables $f : \Omega \rightarrow A$ with a finite alphabet A .
 253 In fact, there is a sequence of random variables f_n with finite alphabets, for example, obtained

254 via increasingly fine partitions of Ω , such that $\text{KL}(P_{f_n} \| M_{f_n})$ tends to $\text{KL}(P \| M)$ as $n \rightarrow \infty$;
 255 see [14, Corollary 5.2.3].

256 *Remark 2.1.* KL is an asymmetric quantity that underlies the definitions of Shannon,
 257 transfer, causation entropy and (conditional) mutual information.

258 A key property is the so-called divergence inequality:

259 **Lemma 2.2 (Divergence inequality, [14] Lemma 5.2.1).** *For any probability measures P, M*
 260 *on a common alphabet, we have $\text{KL}(P \| M) \geq 0$ and the equality holds precisely when $P = M$.*

261 Two cases of KL will be relevant to us.

262 **Lemma 2.3 (Relative entropy density [14] Lemma 5.2.3).** *For any probability measures P, M*
 263 *on a common alphabet, if $P \ll M$, then the Radon-Nikodym derivative $f := dP/dM$ exists,*
 264 *is called the relative entropy density of P with respect to M , and verifies*

$$265 \quad \text{KL}(P \| M) = \int_{\Omega} \ln f(\omega) dP(\omega) = \int_{\Omega} f(\omega) \ln f(\omega) dM(\omega).$$

266 *In this case, if Ω is finite then $f(\omega) = P(\omega)/M(\omega)$ and KL reduces to the finite alphabet case;*
 267 *if $\Omega = \mathbb{R}^d$ and $P, M \ll \text{Leb}$ with densities f_P, f_M , respectively, then*

$$268 \quad \text{KL}(P \| M) = \int_{\mathbb{R}^d} f(x) \ln \frac{f(x)}{g(x)} dx.$$

269 *On the other hand, if P is not absolutely continuous with respect to M , then*

$$270 \quad \text{KL}(P \| M) = \infty.$$

271 **2.2. Mutual information.** Define the *mutual information* between two random variables
 272 X and Y to be

$$273 \quad I(X; Y) := \text{KL}(P_{XY} \| P_X \otimes P_Y).$$

274 It can be shown [14, Chapter 2.5] that the (Shannon) *entropy* of X (defined as $H(X) :=$
 275 $-\sum_{x \in A_X} p_X(x) \ln p_X(x)$ in the discrete alphabet case) can be recovered by the mutual infor-
 276 mation with X itself $H(X) = I(X; X)$ and therefore $I(X; Y) = H(X) + H(Y) - H(X, Y)$.

277 *Remark 2.4.* In light of Lemma 2.2, it is clear that $I(X; Y)$ equals zero precisely when X, Y
 278 are independent and quantifies their deviation from independence otherwise. The product
 279 of marginals $P_X \otimes P_Y$ serves as the reference independent model against which the joint
 280 distribution P_{XY} is compared. More precisely, if (X', Y') has joint distribution $P_X \otimes P_Y$, then
 281 X', Y' are independent and have same marginal distributions as X, Y .

282 **2.3. Conditional mutual information.** First we consider the finite alphabet case: three
 283 random variables X, Y, Z with finite alphabets A_X, A_Y, A_Z , each equipped with the power-set
 284 σ -algebra $\mathcal{B}_{A_*} = 2^{A_*}$, $*$ = X, Y, Z . Define the *conditional mutual information of X, Y given*
 285 Z to be

$$286 \quad (2.1) \quad I(X; Y | Z) := \text{KL}(P_{XYZ} | P_{X \times Y | Z}),$$

287 where $P_{X \times Y|Z}$ is a probability distribution on $A_X \times A_Y \times A_Z$ defined by

$$288 \quad (2.2) \quad P_{X \times Y|Z}(B_X \times B_Y \times B_Z) := \sum_{z \in B_Z} \mathbb{P}(X \in B_X|Z = z)\mathbb{P}(Y \in B_Y|Z = z)\mathbb{P}(Z = z)$$

289 for any $B_X \in \mathcal{B}_{A_X}$, $B_Y \in \mathcal{B}_{A_Y}$ and $B_Z \in \mathcal{B}_{A_Z}$. Here, the conditional probability is the usual
290 one $\mathbb{P}(F|E) = \frac{\mathbb{P}(F \cap E)}{\mathbb{P}(E)}$ provided that $\mathbb{P}(E) > 0$.

291 *Remark 2.5.* As discussed in the Introduction, conditional mutual information is designed
292 to quantify the deviation from conditional independence of X, Y given Z . And $P_{X \times Y|Z}$ is
293 designed to serve as the conditional independent model against which to compare the joint
294 distribution P_{XYZ} , cf. the role of $P_X \otimes P_Y$ in the definition of $I(X; Y)$ as discussed in Remark
295 2.4. More precisely, consider new random variables X', Y', Z' with joint distribution $P_{X \times Y|Z}$
296 and observe

- 297 • X', Y', Z' have the same marginal distributions as X, Y, Z : $P_{X'} = P_X$, $P_{Y'} = P_Y$,
298 $P_{Z'} = P_Z$;
- 299 • X', Y' have the same conditional marginal distributions given Z' as X, Y given Z :
300 $\mathbb{P}(X \in B_X|Z = z) = \mathbb{P}(X' \in B_X|Z' = z)$ and $\mathbb{P}(Y \in B_Y|Z = z) = \mathbb{P}(Y' \in B_Y|Z' = z)$;
- 301 • X', Y' are conditionally independent given Z' : $\mathbb{P}(X' \in B_X, Y' \in B_Y|Z' = z) = \mathbb{P}(X' \in$
302 $B_X|Z' = z)\mathbb{P}(Y' \in B_Y|Z' = z)$.

303 In other words, $P_{X \times Y|Z}$ is a “Markovization” of the joint distribution P_{XYZ} in the sense
304 that the modified random variables X', Y', Z' form a Markov chain $Y' \rightarrow Z' \rightarrow X'$ (or
305 $X' \rightarrow Z' \rightarrow Y'$) because the information about the state of Y' , in addition to that of Z' ,
306 does not further resolve the uncertainty about the state of X' (the same holds with X', Y'
307 swapped).

308 To generalize the definition of $I(X; Y|Z)$ in Eq. 2.1, the main challenge lies with the
309 conditional probabilities appearing in the definition (2.2) of $P_{X \times Y|Z}$. In general, we may well
310 have $\mathbb{P}(Z = z) = 0$ for each $z \in A_Z$, for example, take Z to be uniformly distributed on
311 $A_Z = [0, 1]$, or any other distribution absolutely continuous with respect to Lebesgue. This
312 makes it impossible to define $\mathbb{P}(F|Z = z)$ in the same way as the discrete alphabet case
313 $\frac{\mathbb{P}(F \cap \{Z=z\})}{\mathbb{P}(Z=z)}$.

314 This challenge can be met by (i) interpreting the conditional probability $\mathbb{P}(X \in B_X|Z =$
315 $z)$, rather than a fraction, as a Radon-Nikodym derivative for fixed $B_X \in \mathcal{B}_{A_X}$ and (ii)
316 requiring that the alphabets of X, Y, Z be “standard” measurable spaces so that $\mathbb{P}(X \in$
317 $B_X|Z = z)$ is well-defined as regular conditional probability simultaneously for all $B_X \in \mathcal{B}_{A_X}$
318 and similarly for $\mathbb{P}(Y \in B_Y|Z = z)$. In [14], there is an even more general definition beyond
319 standard alphabets. Since the standard alphabet already covers the practically relevant cases
320 such as Polish spaces, we shall contain our discussion in the standard alphabet case and leave
321 the details in the Appendix.

322 Consider three random variables X, Y, Z on a common probability space $(\Omega, \mathcal{F}, \mathbb{P})$ with
323 standard alphabets (A_X, \mathcal{B}_{A_X}) , (A_Y, \mathcal{B}_{A_Y}) , (A_Z, \mathcal{B}_{A_Z}) , respectively. See Appendix A for de-
324 tails. Define the *conditional average mutual information* as in Eq. (2.1) where the Markoviza-
325 tion $P_{X \times Y|Z}$ is given in terms of regular conditional probabilities, for $B_X \in \mathcal{B}_{A_X}$, $B_Y \in \mathcal{B}_{A_Y}$,

326 $B_Z \in \mathcal{B}_{A_Z}$

$$\begin{aligned}
327 \quad P_{X \times Y|Z}(B_X \times B_Y \times B_Z) &:= \int_{Z^{-1}B_Z} \mathbb{P}(X \in B_X | \sigma(Z)) \mathbb{P}(Y \in B_Y | \sigma(Z)) d\mathbb{P} \\
328 \quad &= \int_{B_Z} \mathbb{P}(X \in B_X | Z = z) \mathbb{P}(Y \in B_Y | Z = z) dP_Z(z).
\end{aligned}$$

329 *Remark 2.6.* Note that $P_{X \times Y|Z}$ is a deterministic probability measure and hence the con-
330 ditional mutual information $I(X; Y|Z)$ is a deterministic object on $A_X \times A_Y \times A_Z$, even though
331 the notation suggests some conditioning. As the construction above shows, the randomness
332 from conditioning on Z is averaged out.

333 In light of Lemma 2.2, $I(X; Y|Z)$ equals zero precisely when X, Y are conditionally inde-
334 pendent given Z and quantifies the deviation from this conditional independence otherwise.

335 Since $P_{XYZ}, P_{X \times Y|Z}$ both have Z -marginals equal to P_Z by construction, they admit
336 disintegrations with respect to P_Z denoted by $(P_{XYZ})_z$ and $(P_{X \times Y|Z})_z$, which coincide with
337 the regular conditional probabilities: for P_Z -a.e. $z \in A_Z$, and all $B_X \in \mathcal{B}_{A_X}, B_Y \in \mathcal{B}_{A_Y}$, we
338 have

$$\begin{aligned}
339 \quad (P_{XYZ})_z(B_X \times B_Y) &= \mathbb{P}(X \in B_X, Y \in B_Y | Z = z), \\
340 \quad (P_{X \times Y|Z})_z(B_X \times B_Y) &= \mathbb{P}(X \in B_X | Z = z) \mathbb{P}(Y \in B_Y | Z = z).
\end{aligned}$$

341 See Appendix A for more details. We will sometimes prefer the disintegration notation to the
342 regular conditional probability notation for clarity of presentation.

343 *Definition 2.7 (Z-conditioned random variables).* For each $z \in A_Z$, define the Z -conditioned
344 random variables X_z, Y_z with alphabets $(A_X, \mathcal{B}_{A_X}), (A_Y, \mathcal{B}_{A_Y})$, respectively, and joint distri-
345 bution

$$346 \quad (2.3) \quad P_{X_z Y_z} := (P_{XYZ})_z, \quad z \in A_Z.$$

347 Then, their marginal distributions are given by

$$\begin{aligned}
348 \quad P_{X_z}(B_X) &= P_{X_z Y_z}(B_X \times A_Y) = \mathbb{P}(X \in B_X | Z = z), \quad z \in A_Z, B_X \in \mathcal{B}_{A_X} \\
349 \quad P_{Y_z}(B_Y) &= P_{X_z Y_z}(A_X \times B_Y) = \mathbb{P}(Y \in B_Y | Z = z), \quad z \in A_Z, B_Y \in \mathcal{B}_{A_Y}.
\end{aligned}$$

350 Hence,

$$351 \quad (P_{X \times Y|Z})_z = P_{X_z} \otimes P_{Y_z}$$

352 and

$$353 \quad I(X_z; Y_z) = \text{KL}(P_{X_z Y_z} \| P_{X_z} \otimes P_{Y_z}) = \text{KL}((P_{XYZ})_z \| (P_{X \times Y|Z})_z).$$

354 The intuition behind the above construction of X_z, Y_z is to consider them as the disintegrated
355 versions of X, Y on the z -slice $A_X \times A_Y \times \{z\}$. The next proposition shows that the conditional
356 mutual information $I(X; Y|Z)$ is the average of $I(X_z; Y_z)$ across all such z -slices.

357 *Proposition 2.8 (Average of disintegrated MI).* Consider three random variables X, Y, Z on
358 a common probability space $(\Omega, \mathcal{F}, \mathbb{P})$ with standard alphabets $(A_X, \mathcal{B}_{A_X}), (A_Y, \mathcal{B}_{A_Y})$, and

359 (A_Z, \mathcal{B}_{A_Z}) , respectively. Then, the conditional mutual information $I(X; Y|Z)$ is the P_Z -
 360 average of mutual information $I(X_z; Y_z) = \text{KL}((P_{XYZ})_z || (P_{X \times Y|Z})_z)$ between the z -conditioned
 361 random variables X_z, Y_z . More precisely, if $(P_{XYZ})_z \ll (P_{X \times Y|Z})_z$ for P_Z -a.e. $z \in A_Z$, then
 362 the Radon-Nikodym derivative

$$363 \quad \frac{d(P_{XYZ})_z}{d(P_{X \times Y|Z})_z}(x, y) = \frac{dP_{XYZ}}{dP_{X \times Y|Z}}(x, y, z) \quad \text{for } P_{X \times Y|Z}\text{-a.e. } (x, y, z)$$

364 and hence

$$365 \quad I(X; Y|Z) = \int_{A_Z} I(X_z; Y_z) dP_Z(z);$$

366 otherwise, there is $B_Z \in \mathcal{B}_{A_Z}$ with $P_Z(B_Z) > 0$ and $I(X_z; Y_z) = \text{KL}((P_{XYZ})_z || (P_{X \times Y|Z})_z) =$
 367 ∞ for each $z \in B_Z$, in which case $I(X; Y|Z) = \infty$.

368 *Proof.* First consider $(P_{XYZ})_z \ll (P_{X \times Y|Z})_z$ for P_Z -a.e. $z \in A_Z$. Then the Radon-
 369 Nikodym derivative $\frac{d(P_{XYZ})_z}{d(P_{X \times Y|Z})_z}$ exists for P_Z -a.e. $z \in A_Z$. Integrating its logarithm against
 370 $(P_{XYZ})_z$ yields, according to Lemma 2.3,

$$371 \quad \text{KL}((P_{XYZ})_z || (P_{X \times Y|Z})_z) = \int_{A_X \times A_Y} \ln \frac{d(P_{XYZ})_z}{d(P_{X \times Y|Z})_z} d(P_{XYZ})_z.$$

372 Further integrating the above equation against P_Z yields, by definition of disintegration,

$$373 \quad \int_{A_Z} \text{KL}((P_{XYZ})_z || (P_{X \times Y|Z})_z) dP_Z(z) = \int_{A_Z} \int_{A_X \times A_Y} \ln \frac{d(P_{XYZ})_z}{d(P_{X \times Y|Z})_z} d(P_{XYZ})_z dP_Z(z)$$

$$374 \quad = \int_{A_X \times A_Y \times A_Z} \ln \frac{d(P_{XYZ})_z}{d(P_{X \times Y|Z})_z} dP_{XYZ}.$$

375 For any $B_X \in \mathcal{B}_{A_X}, B_Y \in \mathcal{B}_{A_Y}, B_Z \in \mathcal{B}_{A_Z}$, by definition of disintegration, we have

$$376 \quad \int_{B_X \times B_Y \times B_Z} \frac{d(P_{XYZ})_z}{d(P_{X \times Y|Z})_z}(x, y) dP_{X \times Y|Z}(x, y, z)$$

$$377 \quad = \int_{B_Z} \int_{B_X \times B_Y} \frac{d(P_{XYZ})_z}{d(P_{X \times Y|Z})_z}(x, y) d(P_{X \times Y|Z})_z(x, y) dP_Z(z)$$

$$378 \quad = \int_{B_Z} (P_{XYZ})_z(B_X \times B_Y) dP_Z(z)$$

$$379 \quad = P_{XYZ}(B_X \times B_Y \times B_Z).$$

380 By uniqueness of Radon-Nikodym derivative, we conclude

$$381 \quad \frac{d(P_{XYZ})_z}{d(P_{X \times Y|Z})_z}(x, y) = \frac{dP_{XYZ}}{dP_{X \times Y|Z}}(x, y, z) \quad \text{for } P_{X \times Y|Z}\text{-a.e. } (x, y, z).$$

382 We continue

$$383 \quad \int_{A_Z} \text{KL}((P_{XYZ})_z || (P_{X \times Y|Z})_z) dP_Z(z) = \int_{A_X \times A_Y \times A_Z} \ln \frac{dP_{XYZ}}{dP_{X \times Y|Z}} dP_{XYZ} = I(X; Y|Z),$$

384 by Lemma 2.3. This proves the first assertion.

385 Now we consider the other case: there is some B_Z with $P_Z(B_Z) > 0$ and for each $z \in B_Z$,
 386 there is some $B_{XY}^{(z)}$ with $(P_{XYZ})_z(B_{XY}^{(z)}) > 0 = (P_{X \times Y|Z})_z$, then the set

$$387 \quad B_{XYZ} := \bigcup_{z \in B_Z} B_{XY}^{(z)} \times \{z\}$$

388 has the property that

$$389 \quad P_{XYZ}(B_{XYZ}) > 0 = P_{X \times Y|Z}(B_{XYZ}).$$

390 In particular, P_{XYZ} is not absolutely continuous with respect to $P_{X \times Y|Z}$. Hence, by Lemma
 391 2.3, we have

$$392 \quad I(X; Y|Z) = \text{KL}(P_{XYZ} \| P_{X \times Y|Z}) = \infty.$$

393 **Example 2.9 (Transfer and causation entropy).** Consider a stochastic process $X = (X_0, X_1, \dots)$ ■
 394 taking values in (A_X, \mathcal{B}_{A_X}) and another stochastic process $Y = (Y_0, Y_1, \dots)$ taking values in
 395 (A_Y, \mathcal{B}_{A_Y}) .

396 As in [5, Chapter 9.8.1], we are interested to quantify the information flow from Y to X
 397 at time t , conditioned on some history $X_t^{(k)} = (X_t, \dots, X_{t-k+1})$ of X itself k -steps into the
 398 past. If there is no such information flow, then $X_{t+1}, Y_t^{(l)}$ should be conditionally independent
 399 given $X_t^{(k)}$, i.e.,

$$400 \quad I(X_{t+1}; Y_t^{(l)} | X_t^{(k)}) = 0.$$

401 Otherwise, the information flow can be quantified by the deviation from conditional indepen-
 402 dence. This motivates our definition

$$403 \quad T_{Y \rightarrow X, t} := I(X_{t+1}; Y_t^{(l)} | X_t^{(k)}) = \text{KL} \left(P_{X_{t+1} X_t^{(k)} Y_t^{(l)}} \parallel P_{X_{t+1} \times Y_t^{(l)} | X_t^{(k)}} \right),$$

404 which has a similar form to the discrete version given by eq. (9.128) in [5]. Variations such
 405 as unlimited memory can also be considered.

406 The causation entropy is a fruitful generalization of TE in the context of a network of
 407 stochastic processes X^v indexed by nodes $v \in V := \{1, \dots, n\}$, where each $X^v = (X_t^v)_{t=0,1,\dots}$.
 408 Given three collections $I, J, K \subseteq V$ of nodes, CE (looking 1 step into the past) [23] is defined
 409 to be

$$410 \quad C_{J \rightarrow I|K, t} := I(X_{t+1}^{(I)}; X_t^{(J)} | X_t^{(K)}) = \text{KL} \left(P_{X_{t+1}^{(I)} X_t^{(K)} X_t^{(J)}} \parallel P_{X_{t+1}^{(I)} \times X_t^{(J)} | X_t^{(K)}} \right).$$

411 In a discovery algorithm, [23] uses $C_{J \rightarrow I|K, t}$ to quantify the information flowing from nodes
 412 J to nodes I conditioned on nodes K , where I nodes are the potential neighbors of J nodes
 413 under consideration and K is the collection of known neighbors of J ; the authors identify the
 414 most likely neighbors of J as the collection I that maximizes $C_{J \rightarrow I|K, t}$.

415 **3. Dynamic determinism.** This section analyzes the conditional mutual information $I(X; Y|Z)$ ■
 416 in the case where $X = T(Y, Z)$ is determined by Y, Z via a measurable function $T : A_Y \times A_Z \rightarrow$
 417 A_X . This is a context particularly relevant to dynamics.

418 **3.1. Mutual information: zero or positive.** Before diving into the conditional mutual
 419 information among three random variables, we first consider two random variables. We begin
 420 with a trivial observation.

421 **Proposition 3.1 (Zero mutual information).** *Let X, Z be two random variables. If $P_X = \delta_{x_0}$
 422 for some $x_0 \in A_X$, then*

$$423 \quad P_{XZ} = P_X \otimes P_Z.$$

424 *In particular, $I(X; Z) = \text{KL}(P_{XZ} \| P_X \otimes P_Z) = 0$.*

425 Note that $P_X = \delta_{x_0}$ is equivalent to $X \equiv x_0$ a.s.; in this case, we may view $X = T(Z)$
 426 for the constant map $T : z \mapsto x_0$. As we will see shortly, this is essentially the only way for
 427 $I(X; Z)$ to vanish.

428 More generally, consider a measurable map $T : A_Z \rightarrow A_X$ and two random variables X, Z .
 429 The following are equivalent

- 430 (i) $X = T(Z)$ a.s.
- 431 (ii) $\mathbb{P}(X = T(Z)|Z) = 1$ a.s.
- 432 (iii) $\mathbb{P}(X = T(Z)|Z = z) = 1$ for P_Z -a.e. $z \in A_Z$.

433 When one of the above holds, we say that X is determined by Z via T .

434 **Proposition 3.2 (Positive mutual information).** *Consider a random variable $X = T(Z)$,
 435 determined by another random variable Z via some measurable map $T : A_Z \rightarrow A_X$. If there
 436 is some $B_X \in \mathcal{B}_{A_X}$ with $0 < P_X(B_X) < 1$, then the event*

$$437 \quad S := B_X \times T^{-1}(A_X \setminus B_X)$$

438 *has the property that*

$$439 \quad P_{XZ}(S) = 0 < P_X \otimes P_Z(S);$$

440 *in particular, we have $P_{XZ} \neq P_X \otimes P_Z$ and $I(X; Z) = \text{KL}(P_{XZ} \| P_X \otimes P_Z) > 0$.*

Proof.

$$441 \quad P_{XZ}(S) = \mathbb{P}((X, Z) \in S) = \mathbb{P}((T(Z), Z) \in B_X \times T^{-1}(A_X \setminus B_X))$$

$$442 \quad = \mathbb{P}(Z \in T^{-1}(B_X) \cap T^{-1}(A_X \setminus B_X)) = 0$$

443 and

$$444 \quad P_X \otimes P_Z(S) = P_X(B_X)P_Z(T^{-1}(A_X \setminus B_X)) = P_X(B_X)P_X(A_X \setminus B_X) > 0.$$

445 This completes the proof. ■

446 Proposition 3.2 provides a partial converse to Proposition 3.1. If we additionally require
 447 that the alphabet (A_X, \mathcal{B}_{A_X}) be such that every zero-one measure is a dirac delta, then it is
 448 a complete converse.

449 A measure μ on a measurable space (A, \mathcal{A}) is said to be a *zero-one measure* if $\mu(F)$ is
 450 either 0 or 1 for all $F \in \mathcal{A}$. A dirac delta is necessarily a zero-one measure, but there are
 451 zero-one measures which are not dirac deltas. The issue usually is that the σ -algebra is too
 452 coarse.

453 *Example 3.3 (Non-measurable singletons).* Consider the alphabet $A_X = \{a, b\}$ equipped
 454 with the trivial σ -algebra $\mathcal{B}_{A_X} = \{\emptyset, A_X\}$. The only probability measure P_X on (A_X, \mathcal{B}_{A_X}) is
 455 a zero-one measure, but not a dirac delta because the singletons $\{a\}, \{b\}$ are not measurable.

456 Combining Propositions 3.1 and 3.2 yields the following dichotomy result.

457 **Theorem 3.4 (Characterization of positive mutual information).** *Let X be a random variable*
 458 *with an alphabet (A_X, \mathcal{B}_{A_X}) , where every zero-one measure is a dirac delta. Suppose also that*
 459 *$X = T(Z)$ is determined by another random variable Z via a measurable map $T : A_Z \rightarrow A_X$.*
 460 *Then, we have a dichotomy:*

- 461 (i) X is constant. In this case, $I(X; Z) = 0$;
 462 (ii) X is nonconstant. In this case, $I(X; Z) > 0$.

463 Discrete spaces and Polish spaces are key examples where every zero-one measure is a
 464 dirac delta.

465 *Example 3.5 (Separable metric space).* If A is a separable metric space, equipped with the
 466 Borel σ -algebra \mathcal{A} , then any zero-one measure on (A, \mathcal{A}) must be a dirac delta. Indeed, if
 467 μ were a zero-one measure on (A, \mathcal{A}) but not a dirac delta, then the support of μ is well-
 468 defined (see [17, Theorem 2.1]) and must contain at least two distinct points $x_1 \neq x_2$ with
 469 $d(x_1, x_2) = d > 0$. By definition of support, the two open balls $B_i := B(x_i, d/3)$ are disjoint
 470 with $\mu(B(x_i, d/3)) = 1 > 0$. Now we arrive at $1 = \mu(A) \geq \mu(B_1) + \mu(B_2) = 1 + 1 = 2$, a
 471 contradiction.

472 Specific examples include a finite or countable set A equipped with the discrete distance
 473 $d(a, b) = \delta_{ab}$, and other Polish spaces equipped with the Borel σ -algebra.

474 3.2. Mutual information: finite or infinite.

475 **Theorem 3.6 (Mutual information: finite or infinite).** *Consider a random variable $X =$*
 476 *$T(Z)$ determined by another random variable Z via some measurable map $T : A_Z \rightarrow A_X$.*
 477 *Assume that the singletons are measurable, i.e., $\{x\} \in \mathcal{B}_{A_X}$ for all $x \in A_X$. Then, we have a*
 478 *dichotomy:*

- 479 1. *Atomic case: there is a finite or countable set $S_X \in \mathcal{B}_{A_X}$ with $P_X(S_X) = 1$. In this*
 480 *case, $P_{XZ} \ll P_X \otimes P_Z$ and*

$$481 \quad I(X; Z) = \sum_{x \in S_X} P_X(x) \int_{A_Z} \ln \frac{d(P_{XZ})_x}{dP_Z}(z) d(P_{XZ})_x(z),$$

482 *which can be either finite or infinite.*

- 483 2. *Continuous case: there is $B_X \in \mathcal{B}_{A_X}$ with $P_X(B_X) > 0$ and $P_X(\{x\}) = 0$ for all*
 484 *$x \in B_X$. In this case, the set*

$$485 \quad S := \{(T(z), z) : T(z) \in B_X\}$$

486 *has the property that $P_{XZ}(S) > 0$ and $P_X \otimes P_Z(S) = 0$. In particular, P_{XZ} is not*
 487 *absolutely continuous with respect to $P_X \otimes P_Z$ and hence*

$$488 \quad I(X; Z) = \text{KL}(P_{XZ} \| P_X \otimes P_Z) = \infty.$$

489 *Proof.* For the atomic case, consider any $N \in \mathcal{B}_{A_X} \otimes \mathcal{B}_{A_Z}$ with $P_X \otimes P_Z(N) = 0$. We
 490 show $P_{XZ}(N) = 0$. By Fubini, for any $x \in S_X$, we have $P_Z(N_x) = 0$, where $N_x := \{z \in A_Z :$
 491 $(x, z) \in N\}$. Hence,

$$492 \quad P_{XZ}(N) = \mathbb{P}((X, Z) \in N) = \sum_{x \in S_X} \mathbb{P}((X, Z) \in \{x\} \times N_x) = \sum_{x \in S_X} \mathbb{P}(T(Z) = x, Z \in N_x)$$

$$493 \quad \leq \sum_{x \in S_X} \mathbb{P}(Z \in N_x) = \sum_{x \in S_X} P_Z(N_x) = 0.$$

494 Now we show that the atomic and continuous cases form a dichotomy. If P_X does not
 495 admit a B_X with $P_X(B_X) > 0$ and $P_X(\{x\}) = 0$ for all $x \in B_X$, then for every B_X with
 496 $P_X(B_X) > 0$, there is some $x \in B_X$ with $P_X(\{x\}) > 0$. Since P_X is a probability measure,
 497 there can be at most countably many $x \in A_X$ with $P_X(\{x\}) > 0$; denote by A_X^1 the set of
 498 all such point atoms x of P_X . Note A_X^1 is measurable because each singleton is measurable.
 499 Then by construction we must have $P_X(A_X \setminus A_X^1) = 0$ because otherwise $A_X \setminus A_X^1$ would
 500 have positive measure and hence contain a point from A_X^1 . This shows that P_X has at most
 501 countable support, namely, A_X^1 , so we are in the atomic case 1. We conclude that the two
 502 cases indeed form a dichotomy.

503 In the atomless case 2, by definition,

$$504 \quad P_{XZ}(S) = \mathbb{P}((X, Z) \in S) = \mathbb{P}(X = T(Z) \in B_X) = P_X(B_X) > 0.$$

505 By Fubini,

$$506 \quad P_X \otimes P_Z(S) = \int_{A_Z} P_X(S_z) dP_Z(z) = \int_{T^{-1}(B_X)} P_X(\{T(z)\}) dP_Z(z) = 0,$$

507 where in the last equality we use $T(z) \in B_X$ for any $z \in T^{-1}(B_X)$. ■

508 **3.3. Infinite conditional mutual information.** In this section, we consider the case when
 509 Y, Z together determine X , that is, $X = T(Y, Z)$ for a measurable map $T : A_Y \times A_Z \rightarrow A_X$.
 510 We split the alphabet A_Z into three disjoint pieces

$$511 \quad A_Z = A_Z^0 \cup A_Z^{\text{atomic}} \cup A_Z^{\text{continuous}},$$

512 where A_Z^0 consists of $z \in A_Z$ for which the marginal distribution $P_{X_z} := \mathbb{P}(X \in \cdot | Z = z)$
 513 of X_z concentrates on a singleton, i.e., $\mathbb{P}(X = x_z | Z = z) = 1$ for some $x_z \in A_X$; A_Z^{atomic}
 514 consists of $z \in A_Z$ for which P_{X_z} concentrates on a non-singleton at most countable set, i.e.,
 515 $P_{X_z}(B_X) = \mathbb{P}(X \in B_X | Z = z) = 1$ for some non-singleton at most countable $B_X \in \mathcal{B}_{A_X}$;
 516 $A_Z^{\text{continuous}}$ consists of $z \in A_Z$ for which P_{X_z} charges an atomless continuum, i.e., there is
 517 $B_X \in \mathcal{B}_{A_X}$ with $P_{X_z}(B_X) > 0$ and $P_{X_z}(\{x\}) = 0$ for all $x \in B_X$. By Theorem 3.6, the three
 518 parts are disjoint and indeed form a partition of A_Z .

519 **Theorem 3.7 (Conditional mutual information).** *Let random variable $X = T(Y, Z)$ be de-*
 520 *termined by random variables Y, Z via a measurable map $T : A_Y \times A_Z \rightarrow A_X$. Suppose X, Y, Z*
 521 *all have standard alphabets. Then,*

$$522 \quad I(X; Y | Z) = \begin{cases} \int_{A_Z^{\text{atomic}}} I(X_z; Y_z) dP_Z(z) & \text{if } P_Z(A_Z^{\text{continuous}}) = 0, \\ \infty & \text{else.} \end{cases}$$

523 In particular, when $P_Z(A_Z^0) = 1$, we have $I(X; Y|Z) = 0$.

524 *Proof.* By Proposition 2.8, we split the conditional mutual information into three parts.

$$525 \quad I(X; Y|Z) = \int_{A_Z^0} I(X_z; Y_z) dP_Z(z) + \int_{A_Z^{\text{atomic}}} I(X_z; Y_z) dP_Z(z) + \int_{A_Z^{\text{continuous}}} I(X_z; Y_z) dP_Z(z)$$

526 By Proposition 3.1, $I(X_z; Y_z) = 0$ for any $z \in A_Z^0$, so the first term vanishes. The last term is
527 zero when $P_Z(A_Z^{\text{continuous}}) = 0$ and is ∞ otherwise, according to Theorem 3.6. ■

528 In many dynamically relevant situations, we have $P_Z(A_Z^{\text{continuous}}) > 0$, as announced in
529 Theorem A, and hence $I(X; Y|Z) = \infty$.

530 4. Examples.

531 **4.1. Bernoulli interval maps.** Consider the piecewise linear expanding map $E_d : [0, 1] \rightarrow$
532 $[0, 1]$, $d \in \mathbb{Z}$, $d \geq 2$, on the unit interval given by $E_d(x) = d \cdot x \pmod{1}$.

533 If Y is uniformly distributed on the interval, i.e., $P_Y = \text{Leb}_{[0,1]}$, then so is $X = E_d(Y)$,
534 i.e., $P_X = P_Y$. The joint distribution of (X, Y) on the unit square $[0, 1]^2$ is given by $P_{XY} =$
535 $(E_d, \text{id})_* \text{Leb}_{[0,1]}$, which is supported on the graph of E_d . In particular, P_{XY} is mutual singular
536 with respect to $P_X \otimes P_Y = \text{Leb}_{[0,1]^2}$. By Theorem 3.6, $I(X; Y) = \text{KL}(P_{XY} \| P_X \otimes P_Y) = \infty$.

537 Now we discretize. Fix a positive integer $L = \Delta^{-1} \in \mathbb{Z}_{>0}$. Then, the uniform partition
538 by $\{[\frac{i-1}{L}, \frac{i}{L}]\}$ is a Markov partition for E_d . Note

$$539 \quad \mathbb{P}(X^\Delta = i\Delta) = \text{Leb}_{[0,1]} \left[\frac{i-1}{L}, \frac{i}{L} \right) = \frac{1}{L}, \quad \forall i = 1, \dots, L.$$

540 This shows that X^Δ is uniformly distributed on $\{1, \dots, L\}$. So is $Y^\Delta = \Pi^\Delta Y$. The joint
541 distribution $P_{X^\Delta Y^\Delta}$ of (X^Δ, Y^Δ) charges uniform mass to the pairs

$$542 \quad (4.1) \quad (d(i-1) + r \pmod{L}, i), \quad i = 1, \dots, L, \quad r = 1, \dots, d.$$

543 When $L \leq d$, then $P_{X^\Delta Y^\Delta}$ is uniform on $\{1, \dots, L\}^2$, with $P_{X^\Delta Y^\Delta}(i, j) = \frac{1}{L^2}$ for each $(j, i) \in$
544 $\{1, \dots, L\}^2$. In this case,

$$545 \quad I(X^\Delta; Y^\Delta) = \text{KL}(P_{X^\Delta Y^\Delta} \| P_{X^\Delta} \otimes P_{Y^\Delta}) = 0.$$

546 When $L > d$, then only dL pairs of $(j, i) \in \{1, \dots, L\}^2$ satisfying eq. (4.1) are charged
547 with mass $\frac{1}{dL}$ each. In this case,

$$548 \quad I(X^\Delta; Y^\Delta) = \text{KL}(P_{X^\Delta Y^\Delta} \| P_{X^\Delta} \otimes P_{Y^\Delta})$$

$$549 \quad = \sum_{(j,i)} P_{X^\Delta Y^\Delta}(j, i) \ln \frac{P_{X^\Delta Y^\Delta}(j, i)}{P_{X^\Delta}(j) P_{Y^\Delta}(i)}$$

$$550 \quad = dL \frac{1}{dL} \ln \frac{1/dL}{1/L^2} = \ln L - \ln d.$$

551 In the discretized version, a more expanding map E_d with large d gives less mutual infor-
552 mation.

553 **4.2. Sine box functions.** Consider the sine box function $S_n : [0, 1] \rightarrow [0, 1]$ given by

$$554 \quad S_n(x) := \frac{1 + \sin 2\pi nx}{2}, \quad n = 1, 2, \dots$$

555 We compute its invariant measure μ_n by taking a long trajectory $\{x_t = S_n^t(x_0) : t = \tau_0, \tau_0 +$
556 $1, \dots, \tau_0 + \tau - 1\}$, starting from $x_0 = 0.5$ (other initial points $0.2, 0.3, \dots, 0.9$ yielded very
557 similar results), discarding the first $\tau_0 = 1000$ iterates as transient, and collecting the next
558 $\tau = 10^6$ iterates to approximate

$$559 \quad \mu_n \approx \mu_n^{(\tau)} := \frac{1}{\tau} \sum_{t=\tau_0}^{\tau_0+\tau-1} \delta_{x_t}.$$

560 If Y follows μ_n , then $X = S_n(Y)$ follows $(S_n)_* \mu_n = \mu_n$.

561 The probability density function ϕ of μ_n is approximated by the histogram for $\{x_t\}$ binned
562 into $\{(i-1)\Delta, i\Delta\}$, that is,

$$563 \quad \phi^{(\tau)}((i-1)\Delta) := \frac{1}{\tau} \sum_{t=\tau_0}^{\tau_0+\tau-1} \mathbb{1}_{[(i-1)\Delta, i\Delta)}(x_t), \quad i = 1, \dots, L,$$

564 which can be represented in vector form

$$565 \quad \phi^{(\tau)} = (\phi_i^{(\tau)})_{i=1}^L, \quad \phi_i^{(\tau)} := \phi^{(\tau)}((i-1)\Delta).$$

566 The product of the marginals $P_X \otimes P_Y$ discretizes into $P_{X\Delta} \otimes P_{Y\Delta} = (\Pi^\Delta \times \Pi^\Delta)(P_X \otimes P_Y)$,
567 which is approximated by

$$568 \quad P_{X\Delta} \otimes P_{Y\Delta} \approx P_{X\Delta}^{(\tau)} \otimes P_{Y\Delta}^{(\tau)} := \phi^{(\tau)} \cdot (\phi^{(\tau)})^\top$$

569 The joint distribution P_{XY} discretizes into $P_{X\Delta Y\Delta} = (\Pi^\Delta, \Pi^\Delta)(P_{XY})$, which is then ap-
570 proximated by

$$571 \quad P_{X\Delta Y\Delta} \approx P_{X\Delta Y\Delta}^{(\tau)} = (P_{X\Delta Y\Delta}^{(\tau)})_{i,j=1}^L,$$

572 where

$$573 \quad (P_{X\Delta Y\Delta}^{(\tau)})_{i,j} = \frac{1}{\tau} \sum_{t=\tau_0}^{\tau_0+\tau-1} \mathbb{1}_{[(i-1)\Delta, i\Delta)}(x_t) \cdot \mathbb{1}_{[(j-1)\Delta, j\Delta)}(x_{t+1}).$$

574 It follows from Theorem 3.6 that $I(X, Y) = \infty$ for any n . However, higher value of n decreases
575 the ability to resolve uncertainty about $X = S_n(Y)$ from knowledge about Y . Accordingly,
576 we expect $I(X^\Delta; Y^\Delta)$ to decrease as n increases. This is confirmed by simulations as shown
577 in Figure 3.

578 **Appendix A. Regular conditional probability and disintegration on standard measurable**
 579 **spaces.** We motivate the consideration of standard measurable spaces by an attempt to gen-
 580 eralize the definition of conditional probability for discrete variables to more general variables.
 581 We finish the discussion by showing that regular conditional probabilities are equivalent to
 582 disintegrations in our setting.

583 Consider a common probability space $(\Omega, \mathcal{F}, \mathbb{P})$ for random variables X, Y, Z taking values
 584 in $(A_X, \mathcal{B}_{A_X}), (A_Y, \mathcal{B}_{A_Y}), (A_Z, \mathcal{B}_{A_Z})$.

585 The first challenge in generalizing the definition of conditional probability $\mathbb{P}(F|Z = z) :=$
 586 $\frac{\mathbb{P}(F \cap \{Z = z\})}{\mathbb{P}(Z = z)}$ to non-discrete variables Z is that the events $\{Z = z\}$ being conditioned on may
 587 well have zero probability. To overcome this challenge, a first fix is to interpret the conditional
 588 probability as a density (Radon-Nikodym derivative) rather than a fraction. More precisely,
 589 given an arbitrary random variable Z and a fixed event $F \in \mathcal{F}$, we define $\mathbb{P}(F|Z = z), z \in A_Z$
 590 to be the Radon-Nikodym derivative

$$591 \quad \mathbb{P}(F|Z = z) := \frac{d\mathbb{P}^F(Z \in \cdot)}{dP_Z}(z) = \frac{d\mathbb{P}^F(Z \in \cdot)}{dP_Z}(z), \quad z \in A_Z,$$

592 where $\mathbb{P}^F(Z \in \cdot) := \mathbb{P}(F \cap \{Z \in \cdot\})$ is absolutely continuous with respect to $P_Z = \mathbb{P}(Z \in \cdot)$.
 593 By Radon-Nikodym Theorem, $\mathbb{P}(F|Z = z)$ exists and is P_Z -essentially unique. Equivalently,
 594 we have the defining equation for $\mathbb{P}(F|Z = z)$

$$595 \quad \mathbb{P}(F \cap \{Z \in B_Z\}) = \mathbb{P}^F(Z \in B_Z) = \int_{B_Z} \mathbb{P}(F|Z = z) dP_Z(z), \quad B_Z \in \mathcal{B}_{A_Z},$$

596 an analogue of the discrete alphabet case

$$597 \quad \mathbb{P}(F \cap \{Z \in B_Z\}) = \sum_{z \in B_Z} \mathbb{P}(F|Z = z) P_Z(z), \quad B_Z \in \mathcal{B}_{A_Z}.$$

598 This is a more direct construction than the usual conditioning on sigma-algebra, which
 599 we review below for comparison. For a fixed event $F \in \mathcal{F}$, the *conditional probability* $\mathbb{P}(F|\mathcal{G})$
 600 given a sigma-algebra $\mathcal{G} \subseteq \mathcal{F}$ is defined to be any \mathcal{G} -measurable random variable $g : \Omega \rightarrow [0, 1]$
 601 with

$$602 \quad \int_G g d\mathbb{P} = \mathbb{P}(F \cap G), \quad \forall G \in \mathcal{G}.$$

603 $\mathbb{P}(F|\mathcal{G})$ exists and is \mathbb{P} -a.s. unique as the Radon-Nikodym derivative of $\mathbb{P}(F \cap \cdot)/\mathbb{P}(F)$ with
 604 respect to \mathbb{P} , both restricted to \mathcal{G} , provided $\mathbb{P}(F) > 0$; in case $\mathbb{P}(F) = 0$, we have $\mathbb{P}(F|\mathcal{G}) \equiv 0$.
 605 Now consider $\mathcal{G} = \sigma(Z)$. Since $\mathbb{P}(F|\sigma(Z))$ is $\sigma(Z)$ -measurable, it can be factored through Z
 606 [13, Lemma 5.2.1], that is,

$$607 \quad \mathbb{P}(F|\sigma(Z)) = h \circ Z,$$

608 for some measurable function $h : A_Z \rightarrow [0, 1]$. We thus have

$$609 \quad \mathbb{P}(F|Z = z) = h(z).$$

610 A subtle issue remains with this Radon-Nikodym construction, namely, the potential pile
 611 up of exceptional sets $E(F)$ in the definition of $\mathbb{P}(F|Z = z)$. The Radon-Nikodym derivative

612 $\mathbb{P}(F|Z = z)$ is well-defined up to an exceptional set $E(F)$ with $\mathbb{P}(E(F)) = 0$ depending on the
 613 event F . These exceptional sets may pile up $\mathbb{P}(\bigcup_{F \in \mathcal{F}} E(F)) = 1$ and in this case we cannot
 614 define $\mathbb{P}(F|Z = z)$ simultaneously for all $F \in \mathcal{F}$. An example of such a pathology can be
 615 found in [9, Page 624]; for more details see [10, Chapter 5.1.3]. Hence, in order to generalize
 616 the definition of $P_{X \times Y|Z}$ as in Eq. (2.2), we need to rule out such pathologies. This motivates
 617 our second fix: the regular conditional probability.

618 **Definition A.1 (Regular conditional probability (RCP); [13] Chapter 5.8).** The *regular con-*
 619 *ditional probability* given a sub- σ -algebra $\mathcal{G} \subseteq \mathcal{F}$ is a function $f : \mathcal{F} \times \Omega \rightarrow [0, 1]$ such that

- 620 1. for each $\omega \in \Omega$, $f(\cdot, \omega)$ is a probability measure on (Ω, \mathcal{F}) ;
- 621 2. for each $F \in \mathcal{F}$, $f(F, \cdot)$ is a version of $\mathbb{P}(F|\mathcal{G})$.

622 We consider sigma-algebra $\mathcal{G} = \sigma(Z)$ and events of the form $F = \{X \in B_X\} \in \mathcal{F}$. Define the
 623 *regular conditional distribution* of X given Z to be

$$624 \quad \mathbb{P}(X \in B_X|Z = z) := f(\{X \in B_X\}, \omega), \quad \omega \in Y^{-1}\{z\}.$$

625 RCP does not always exist in general but it does, for example, [13, Corollary 5.8.1] (i)
 626 when both (A_X, \mathcal{B}_{A_X}) and (A_Z, \mathcal{B}_{A_Z}) are standard, (ii) when either is discrete.

627 **Definition A.2 (Standard measurable space; [2] page 541).** A measurable space (Ω, \mathcal{F}) is
 628 called a *standard measurable space* if isomorphic via a bi-measurable bijection to a Borel subset
 629 of a Polish space.

630 In particular, a standard measurable space (Ω, \mathcal{F}) admits a sequence of finite fields $\mathcal{F}_n \subseteq \mathcal{F}$,
 631 $n = 0, 1, \dots$ such that

- 632 1. increasing fields: $\mathcal{F}_n \subseteq \mathcal{F}_{n+1}$ for all $n = 0, 1, \dots$;
- 633 2. generating fields: $\mathcal{F} = \sigma(\bigcup_{n=0}^{\infty} \mathcal{F}_n)$;
- 634 3. nonempty atomic intersection: an event is called an *atom* of a field if it is nonempty
 635 and its only subsets which are members of the field are the empty set and itself. If
 636 $G_n \in \mathcal{F}_n$, $n = 0, 1, \dots$ are atoms with $G_{n+1} \subseteq G_n$ for all n , then

$$637 \quad \bigcap_{n=0}^{\infty} G_n \neq \emptyset.$$

638 In fact, the above three conditions are sometimes taken to be the defining properties of a
 639 standard measurable space, for example in [14]. We have taken the more restricted definition
 640 of Arnold [2] to ensure that both regular conditional probabilities and disintegrations exist.

641 Now we review disintegrations and show that they coincide with regular conditional prob-
 642 abilities in our setting.

643 **Definition A.3 (Disintegration; [2] pp 22).** Given a probability measure μ on a product
 644 measurable space $(A \times B, \mathcal{A} \otimes \mathcal{B})$ and a probability measure ν on (A, \mathcal{A}) , we say that a
 645 function $\mu_a(\cdot) : B \rightarrow [0, 1]$ is a *disintegration* of μ with respect to ν if

- 646 1. for all $B \in \mathcal{B}$, $a \mapsto \mu_a(B)$ is measurable function from (A, \mathcal{A}) to $([0, 1], \mathcal{B}([0, 1]))$;
- 647 2. for ν -a.e. $a \in A$, $B \mapsto \mu_a(B)$ is a probability measure on (B, \mathcal{B}) ;
- 648 3. for all $E \in \mathcal{A} \otimes \mathcal{B}$,

$$649 \quad \mu(E) = \int_A \int_B \mathbb{1}_E(a, b) d\mu_a(b) d\nu(a).$$

650 Disintegrations do not always exist, but they do exist ν -essentially uniquely, when (A, \mathcal{A}) ,
 651 (B, \mathcal{B}) are both standard alphabets, see [2, Proposition 1.4.3] and [13, Corollary 5.8.1].

652 Returning to our previous setting, $P_{XYZ}, P_{X \times Y|Z}$ both have Z -marginals equal to P_Z by
 653 construction and so both admit disintegrations with respect to P_Z denoted by $(P_{XYZ})_z$ and
 654 $(P_{X \times Y|Z})_z$. In this case, it follows from the definitions of RCP and disintegration and their
 655 existence and essential uniqueness that for P_Z -a.e. $z \in A_Z$, and all $B_X \in \mathcal{B}_{A_X}, B_Y \in \mathcal{B}_{A_Y}$, we
 656 have

$$\begin{aligned} 657 & (P_{XYZ})_z(B_X \times B_Y) = \mathbb{P}(X \in B_X, Y \in B_Y | Z = z), \\ 658 & (P_{X \times Y|Z})_z(B_X \times B_Y) = \mathbb{P}(X \in B_X | Z = z) \mathbb{P}(Y \in B_Y | Z = z). \end{aligned}$$

659 **Appendix B. Additive noise.** Consider a measurable map $T_0 : [0, 1] \rightarrow [0, 1]$ on the
 660 unit interval, which is *nonsingular* with respect to the Lebesgue measure λ on $[0, 1]$ in the
 661 sense that $\lambda(T_0^{-1}N) = 0$ for any $\lambda(N) = 0$. Consider random variable Z with distribution
 662 $P_Z = h_Z \lambda$.

663 Let $X_0 = T_0(Z)$. By Theorem 3.6, we have $I(X_0; Z) = \infty$.

664 Now perturb T_0 by additive noise

$$665 \quad T_\xi : z \mapsto T_0(z) + \xi \pmod{1},$$

666 where the noise ξ is independent of Z and follows some distribution $P_\xi = h_\xi \lambda$.

667 For concreteness, we take the uniform noise of amplitude ϵ centered at 0 with density
 668 $h_\xi = \frac{1}{\epsilon} \mathbb{1}_{[-\epsilon/2, \epsilon/2]}$.

669 Consider X given by the randomly transformed Z via $\{T_\xi\}$; more precisely,

$$670 \quad \mathbb{P}(X \in B | Z = z) = \int_0^1 \mathbb{1}_B \circ T_\xi(z) dP_\xi(\xi).$$

671 In other words,

$$672 \quad (P_{XZ})_z = (R_{T_0(z)})_* P_\xi, \quad R_\alpha : x \mapsto x + \alpha \pmod{1}.$$

673 If the joint distribution $P_{XZ} \ll P_X \otimes P_Z$, then

$$674 \quad I(X; Z) = \int_{[0,1]^2} f \ln f dP_X \otimes P_Z,$$

675 where $f(x, z) = \frac{dP_{XZ}}{dP_X \otimes P_Z}(x, z) = \frac{d(P_{XZ})_z}{d(P_X \otimes P_Z)_z}(x) = \frac{d(\frac{1}{\epsilon} \mathbb{1}_{[T_0(z) - \epsilon/2, T_0(z) + \epsilon/2]}) \lambda}{dP_X}(x)$.

676 In general, f depends on T_0 . Consider the special case of Bernoulli maps $T_0 = E_d$ or
 677 rotations $T_0 = R_\alpha$, both of which preserve λ . Then, $P_X = \lambda$, $f(x, z) = \frac{1}{\epsilon} \mathbb{1}_{[T_0(z) - \epsilon/2, T_0(z) + \epsilon/2]}(x)$,
 678 and we have

$$\begin{aligned} 679 \quad I(X; Z) &= \int_{[0,1]^2} \frac{1}{\epsilon} \mathbb{1}_{[T_0(z) - \epsilon/2, T_0(z) + \epsilon/2]}(x) \ln \frac{1}{\epsilon} \mathbb{1}_{[T_0(z) - \epsilon/2, T_0(z) + \epsilon/2]}(x) dx dz \\ 680 &= \int_0^1 \int_{T_0(z) - \epsilon/2}^{T_0(z) + \epsilon/2} \frac{1}{\epsilon} \ln \frac{1}{\epsilon} dx dz \\ 681 &= \epsilon \frac{1}{\epsilon} \ln \frac{1}{\epsilon} = \ln \frac{1}{\epsilon}. \end{aligned}$$

682 This indicates that the mutual information of the blurred variables does not distinguish be-
683 tween very ambiguous map $T_0 = E_d$ and non-ambiguous map $T_0 = R_\alpha$.

684 **Appendix C. Derivation of the discretized mutual information formula.**

685 Recall that the Shannon entropy of a continuous random variable X is infinite, but there
686 is a meaningful notion of differential entropy, which differs from the Shannon entropy of the
687 discretization of X by an infinite offset.

688 In a similar spirit, we aim to identify such an infinite offset in mutual information $I(X; Y)$
689 with $X = T(Y)$ so as to extract the meaningful term $A_T(Y)$, which we have termed the
690 relative ambiguity of the system (T, Y) .

691 Observe that $P_{X\Delta Y\Delta} \ll P_{X\Delta} \otimes P_{Y\Delta}$ and hence

$$692 \text{ (C.1) } \quad I(X^\Delta; Y^\Delta) = \sum_{i,j} \mathbb{P}(X^\Delta = i\Delta, Y^\Delta = j\Delta) \ln \frac{\mathbb{P}(X^\Delta = i\Delta, Y^\Delta = j\Delta)}{\mathbb{P}(X^\Delta = i\Delta)\mathbb{P}(Y^\Delta = j\Delta)}.$$

693 Since the densities f_X, f_Y are continuous by assumption in Conjecture C, we have the usual
694 Riemman sum approximation

$$695 \quad \mathbb{P}(X^\Delta = i\Delta) \approx f_X(i\Delta)\Delta$$

$$696 \quad \mathbb{P}(Y^\Delta = j\Delta) \approx f_Y(j\Delta)\Delta.$$

697 In the linear case $T = E_d$, the mass $f_X(i\Delta)\Delta$ splits evenly into $d = |T'(i\Delta)|$ pieces. Since
698 T is piecewise C^1 expanding $|T'| \geq 1$ by assumption in Conjecture C, we conjecture the key
699 approximation

$$700 \quad \mathbb{P}(X^\Delta = i\Delta, Y^\Delta = j\Delta) \approx \frac{f_X(i\Delta)\Delta}{|T'(i\Delta)|}, \quad T(i\Delta) \approx j\Delta.$$

701 When T has contracting regions $|T'| < 1$, this approximation fails. This suggests a connection
702 to the transfer operator formula for expanding maps

$$703 \quad (\hat{T}f)(y) = \sum_{x \in T^{-1}y} \frac{f(x)}{|T'(x)|},$$

704 where the transfer operator $\hat{T} : L^1(\lambda) \rightarrow L^1(\lambda)$ is defined to be the Radon-Nikodym derivative

$$705 \quad \hat{T}f := \frac{dT_*(f\lambda)}{d\lambda}.$$

706 Now we combine these approximations together:

$$\begin{aligned}
707 \quad I(X^\Delta, Y^\Delta) &= \sum_{i\Delta, j\Delta} \mathbb{P}(X^\Delta = i\Delta, Y^\Delta = j\Delta) \ln \frac{\mathbb{P}(X^\Delta = i\Delta, Y^\Delta = j\Delta)}{\mathbb{P}(X^\Delta = i\Delta)\mathbb{P}(Y^\Delta = j\Delta)} \\
708 \quad &\approx \sum_{j\Delta} \sum_{i\Delta \in T^{-1}j\Delta} \frac{f_X(i\Delta)\Delta}{|T'(i\Delta)|} \ln \frac{f_X(i\Delta)\Delta/|T'(i\Delta)|}{f_X(i\Delta)\Delta f_Y(j\Delta)\Delta} \\
709 \quad &= \sum_{j\Delta} \sum_{i\Delta \in T^{-1}j\Delta} \frac{f_X(i\Delta)\Delta}{|T'(i\Delta)|} \ln \frac{1}{|T'(i\Delta)|f_Y(j\Delta)\Delta} \\
710 \quad &\approx \int_{A_Y} \hat{T} \left[f_X \ln \frac{1}{|T'| \cdot f_Y \circ T \cdot \Delta} \right] dy \\
711 \quad &= \int_{A_X} f_X \ln \frac{1}{|T'|} dx + \int_{A_X} f_X \ln \frac{1}{f_Y \circ T} dx + \int_{A_X} f_X \ln \frac{1}{\Delta} dx \\
712 \quad &= - \int_{A_X} f_X \ln |T'| dx + \int_{A_Y} (\hat{T} f_X) \ln \frac{1}{f_Y} dy + \ln \Delta^{-1} \\
713 \quad &= H(Y) - \int \ln |T'| dP_X + \ln \Delta^{-1}.
\end{aligned}$$

714

REFERENCES

- 715 [1] A. A. R. ALMOMANI, J. SUN, AND E. BOLLT, *How entropic regression beats the outliers problem in*
716 *nonlinear system identification*, *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 30 (2020).
717 [2] L. ARNOLD, *Random Dynamical Systems*, Springer Berlin Heidelberg, 1998, [https://doi.org/10.1007/](https://doi.org/10.1007/978-3-662-12878-7)
718 [978-3-662-12878-7](https://doi.org/10.1007/978-3-662-12878-7), <http://dx.doi.org/10.1007/978-3-662-12878-7>.
719 [3] A. ASSAF, M. H. BILGIN, AND E. DEMIR, *Using transfer entropy to measure information flows between*
720 *cryptocurrencies*, *Physica A: Statistical Mechanics and its Applications*, 586 (2022), p. 126484, [https://](https://doi.org/10.1016/j.physa.2021.126484)
721 doi.org/10.1016/j.physa.2021.126484, <http://dx.doi.org/10.1016/j.physa.2021.126484>.
722 [4] L. BARNETT, A. B. BARRETT, AND A. K. SETH, *Granger causality and transfer entropy are equivalent*
723 *for gaussian variables*, *Physical review letters*, 103 (2009), p. 238701.
724 [5] E. M. BOLLT AND N. SANTITISSADEEKORN, *Applied and Computational Measurable Dynamics*, Society
725 for Industrial and Applied Mathematics, Nov. 2013, <https://doi.org/10.1137/1.9781611972641>, [http://](http://dx.doi.org/10.1137/1.9781611972641)
726 dx.doi.org/10.1137/1.9781611972641.
727 [6] T. BOSSOMAIER, L. BARNETT, M. HARRÉ, AND J. T. LIZIER, *An Introduction to Transfer Entropy*,
728 Springer International Publishing, 2016, <https://doi.org/10.1007/978-3-319-43222-9>, [http://dx.doi.](http://dx.doi.org/10.1007/978-3-319-43222-9)
729 [org/10.1007/978-3-319-43222-9](http://dx.doi.org/10.1007/978-3-319-43222-9).
730 [7] T. M. COVER AND J. A. THOMAS, *Elements of Information Theory*, Wiley, Apr. 2005, [https://doi.org/](https://doi.org/10.1002/047174882x)
731 [10.1002/047174882x](https://doi.org/10.1002/047174882x), <http://dx.doi.org/10.1002/047174882x>.
732 [8] T. DIMPFL AND F. J. PETER, *Using transfer entropy to measure information flows between finan-*
733 *cial markets*, *Studies in Nonlinear Dynamics and Econometrics*, 17 (2013), [https://doi.org/10.1515/](https://doi.org/10.1515/snde-2012-0044)
734 [snde-2012-0044](https://doi.org/10.1515/snde-2012-0044), <http://dx.doi.org/10.1515/snde-2012-0044>.
735 [9] J. L. DOOB, *Stochastic Processes*, Wiley Classics Library, John Wiley & Sons, Nashville, TN, Jan. 1990.
736 [10] R. DURRETT, *Cambridge series in statistical and probabilistic mathematics: Probability: Theory and*
737 *examples*, Cambridge University Press, Cambridge, England, 4 ed., Aug. 2010.
738 [11] C. W. GRANGER, *Investigating causal relations by econometric models and cross-spectral methods*, *Econo-*
739 *metrica: journal of the Econometric Society*, (1969), pp. 424–438.
740 [12] C. W. GRANGER, *Some recent development in a concept of causality*, *Journal of econometrics*, 39 (1988),
741 pp. 199–211.

- 742 [13] R. M. GRAY, *Probability, Random Processes, and Ergodic Properties*, Springer US, 2009, <https://doi.org/10.1007/978-1-4419-1090-5>, <http://dx.doi.org/10.1007/978-1-4419-1090-5>.
- 743
- 744 [14] R. M. GRAY, *Entropy and Information Theory*, Springer US, 2011, <https://doi.org/10.1007/978-1-4419-7970-4>, <http://dx.doi.org/10.1007/978-1-4419-7970-4>.
- 745
- 746 [15] D. F. HENDRY, *The nobel memorial prize for clive wj granger*, *Scandinavian Journal of Economics*, 106
747 (2004), pp. 187–213.
- 748 [16] W. M. LORD, J. SUN, N. T. OUELLETTE, AND E. M. BOLLT, *Inference of causal information flow in*
749 *collective animal behavior*, *IEEE Transactions on Molecular, Biological, and Multi-Scale Communi-*
750 *cations*, 2 (2016), pp. 107–116.
- 751 [17] K. PARTHASARATHY, *Probability Measures on Metric Spaces*, Academic Press, 1967.
- 752 [18] Y. B. PESIN, *Characteristic lyapunov exponents and smooth ergodic theory*, *Russian Mathematical Sur-*
753 *veys*, 32 (1977), p. 55–114, <https://doi.org/10.1070/rm1977v032n04abeh001639>, <http://dx.doi.org/10.1070/RM1977v032n04ABEH001639>.
- 754
- 755 [19] T. SCHREIBER, *Measuring information transfer*, *Physical Review Letters*, 85 (2000), p. 461–464, <https://doi.org/10.1103/physrevlett.85.461>, <http://dx.doi.org/10.1103/PhysRevLett.85.461>.
- 756
- 757 [20] D. P. SHORTEN, R. E. SPINNEY, AND J. T. LIZIER, *Estimating transfer entropy in continuous time*
758 *between neural spike trains or other event-based data*, *PLOS Computational Biology*, 17 (2021),
759 p. e1008054, <https://doi.org/10.1371/journal.pcbi.1008054>, <http://dx.doi.org/10.1371/journal.pcbi.1008054>.
- 760
- 761 [21] J. SUN AND E. M. BOLLT, *Causation entropy identifies indirect influences, dominance of neighbors and*
762 *anticipatory couplings*, *Physica D: Nonlinear Phenomena*, 267 (2014), pp. 49–57.
- 763 [22] J. SUN, C. CAFARO, AND E. M. BOLLT, *Identifying the coupling structure in complex systems through*
764 *the optimal causation entropy principle*, *Entropy*, 16 (2014), pp. 3416–3433.
- 765 [23] J. SUN, D. TAYLOR, AND E. M. BOLLT, *Causal network inference by optimal causation entropy*, *SIAM*
766 *Journal on Applied Dynamical Systems*, 14 (2015), p. 73–106, <https://doi.org/10.1137/140956166>,
767 <http://dx.doi.org/10.1137/140956166>.
- 768 [24] S. SURASINGHE AND E. M. BOLLT, *On geometry of information flow for causal inference*, *Entropy*, 22
769 (2020), p. 396, <https://doi.org/10.3390/e22040396>, <http://dx.doi.org/10.3390/e22040396>.
- 770 [25] M. URSINO, G. RICCI, AND E. MAGOSSO, *Transfer entropy as a measure of brain connectivity: A critical*
771 *analysis with the help of neural mass models*, *Frontiers in Computational Neuroscience*, 14 (2020),
772 <https://doi.org/10.3389/fncom.2020.00045>, <http://dx.doi.org/10.3389/fncom.2020.00045>.
- 773 [26] R. VICENTE, M. WIBRAL, M. LINDNER, AND G. PIPA, *Transfer entropy—a model-free measure of effective*
774 *connectivity for the neurosciences*, *J. Comput. Neurosci.*, 30 (2011), pp. 45–67.