

# CAUSAL NETWORK INFERENCE BY OPTIMAL CAUSATION ENTROPY\*

JIE SUN<sup>†</sup>, DANE TAYLOR<sup>‡</sup>, AND ERIK M. BOLLT<sup>§</sup>

**Abstract.** The broad abundance of time series data, which is in sharp contrast to limited knowledge of the underlying network dynamic processes that produce such observations, calls for a rigorous and efficient method of causal network inference. Here we develop mathematical theory of causation entropy, an information-theoretic statistic designed for model-free causality inference. For stationary Markov processes, we prove that for a given node in the network, its causal parents forms the *minimal set of nodes that maximizes causation entropy*, a result we refer to as the *optimal causation entropy principle*. Furthermore, this principle guides us to develop computational and data efficient algorithms for causal network inference based on a two-step discovery and removal algorithm for time series data for a network-couple dynamical system. Validation in terms of analytical and numerical results for Gaussian processes on large random networks highlight that inference by our algorithm outperforms previous leading methods including conditioned Granger causality and transfer entropy. Interestingly, our numerical results suggest that the number of samples required for accurate inference depends strongly on network characteristics such as the density of links and information diffusion rate and not necessarily on the number of nodes.

**Key words.** causal network inference, optimal causation entropy, stochastic network dynamics

**AMS subject classifications.** 37N99, 62B10, 94A17

**1. Introduction.** Research of dynamic processes on large-scale complex networks has attracted considerable interest in recent years with exciting developments in a wide range of disciplines in social, scientific, engineering, and medical fields [48, 49, 74]. One important line of research focuses on exploring the role of network structure in determining the dynamic properties of a system [6, 17, 18, 19, 27, 55, 67, 79] and utilizing such knowledge in controlling network dynamics [15, 70] and optimizing network performance [13, 38, 50, 56, 72]. In applications such as the study of neuronal connectivity or gene interactions, it is nearly impossible to directly identify the network structure without severely interfering with the underlying system whereas time series measurements of the individual node states are often more accessible [68]. From this perspective, it is crucial to reliably infer the network structure that shapes the dynamics of a system from time series data. It is essential that one accounts for directed “cause and effect” relationships, which often offer deeper insight than non-directed relationships (e.g., correlations) [53, 62, 66]. In particular, causal network inference is considered a central problem in the research of social perception [35], epidemiological factors [57], neural connectivity [11, 12], economic impacts [34], and basic physical relationships of climatological events [60, 61]. Evidently, understanding causality is a necessary and important precursor step towards the goal of effectively controlling and optimizing system dynamics (e.g., medical intervention of biological processes and policy design for economic growth and social development).

In a network dynamic process involving a large number of nodes, causal relationships are inherently difficult to infer. For example, the fact that a single node can potentially be influenced by many (if not all) others through network interactions

---

\*This work was funded by ARO Grant No. 61386-EG (J.S and E.M.B), and NSF Grant No. DMS-1127914 through the Statistical and Applied Mathematical Sciences Institute (D.T.).

<sup>†</sup>Department of Mathematics, Clarkson University, Potsdam, NY 13699 ([sunj@clarkson.edu](mailto:sunj@clarkson.edu)).

<sup>‡</sup>Statistical and Applied Mathematical Sciences Institute, Research Triangle Park, NC 27709; Department of Mathematics, University of North Carolina, Chapel Hill, NC 27599.

<sup>§</sup>Department of Mathematics, Clarkson University, Potsdam, NY 13699.

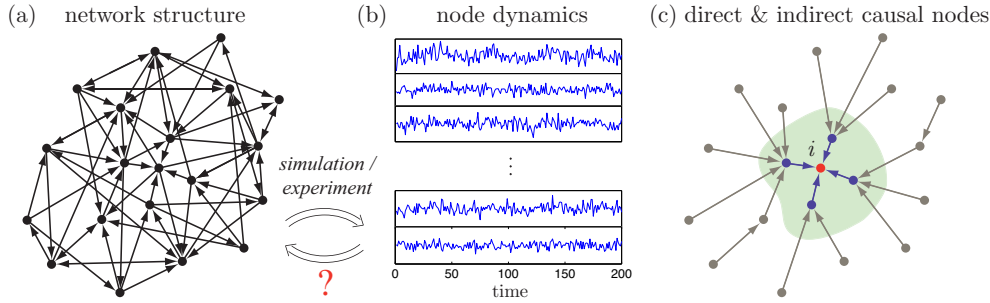


FIG. 1.1. *Network dynamics, time series, and the causal network inference problem.* Modern scientific approaches such as simulation, experiments, and data mining have produced an abundance of high-dimensional time series data describing dynamic processes on complex networks (a→b). Given empirical observations, an important problem is to infer the causal network structure that underlies the observed time series. As shown in (c), for each node  $i$ , the goal is to identify its “causal parents”, the nodes that directly influence its dynamics (nodes in shaded region), while pruning away the nodes that do not (nodes outside the shaded region), thus recovering the direct links to node  $i$  in the causal network. The key to efficiently and accurately identify direct causal links from non-causal ones is to follow an algorithm involving tests for independence via judiciously selected conditioning sets. The main goal of this paper is to develop and validate such algorithms for stationary Markov processes.

makes it challenging to untangle the direct causal links from indirect and erroneous ones (see Fig. 1.1 for illustration). Granger recognized the crucial role played by conditioning and defines a causal relationship based on two basic principles [29, 30]:

- (i) The cause should occur before the effect;
- (ii) The cause should contain information about the caused that is not available otherwise.

A relationship that fulfills both requirements is unambiguously defined as causal. In practice, although the first requirement is straightforward to examine when temporal ordering of the data is available, it is difficult to check the second as it involves the consideration of *all* available information (time series data from all variables). Trade-offs are often made, by either restricting to small-scale networks with no time delay and just a handful of variables [32, 71], or partially removing the second requirement therefore reducing the accuracy of network inference [76]. Inferring large-scale networks from time series data remains to be a relatively open problem [41, 68].

The classical Granger causality test was designed for linear regression models [29, 30], although several extensions have been proposed to nonlinear models, including local linear approximations [14] and partial functional expansion via radial basis functions [2]. Information-based causality inference measures represent a systematically way of overcoming the model-dependent limitation in the linear Granger causality test. In particular, Schreiber proposed transfer entropy as a measure of information flow, or effective coupling, between two processes regardless of the actual functional relationship between them [63]. The transfer entropy from process  $Y$  to  $X$  measures the uncertainty reduction of the future states of  $X$  as a result of knowing the past of  $Y$  given that the past of  $X$  is already known, and is essentially the mutual information between the future of  $X$  and history of  $Y$  *conditioning* on the history of  $X$  [37, 51]. Because of its ability to associate temporal and spatial directionality with coupling, transfer entropy has quickly started to gain popularity in a broad range of disciplines including bioinformatics, neuroscience, climatology and others, as a tool to infer ef-

fective pairwise coupling that underlie complex dynamic processes [9, 76]. However, transfer entropy, which was introduced specifically for detecting the directionality of information flow between two processes, has fundamental limitations when applied in a multivariate setting, to the inference of networks [65, 71]. In particular, without proper conditioning, inference based on transfer entropy tends to produce systematic errors due to, for example, the effects of indirect influences and dominance of neighbors [71]. As shown in Fig. 1.1(c), the main purpose of this work is to identify for each node  $i$  its “causal parents” that directly influence node  $i$ , while not falsely inferring indirect (i.e., non-causal) nodes.

Proper conditioning can distinguish between direct and indirect causal relationships, and it is thus unsurprising that conditioning is widely adopted as a key ingredient in many network inference methods [21, 32, 37, 51, 60, 61, 65, 66, 71]; however, even within such a general theme, the inference of networks requires a theoretically sound approach that is also algorithmically reliable and efficient. For example, one must develop a strategy for choosing which potential links to examine and which nodes to condition on. Thus we note two essential steps in causal network inference: (1) adopting a statistic for the inference of a causal relationship, and (2) developing an algorithm that iteratively employs step (1) to learn the causal network. Whereas accuracy, tractability, and generality of the chosen statistic is often the priority for (1), various challenges arise regarding (2). In particular, these often include minimizing the computational cost by reducing the number of statistics that needs to be computed, as well as reducing the error incurred by finite-sized data by keeping the size of the conditioning sets (i.e., the dimension of the estimation problem) as small as possible. In general, the inaccuracy when estimating statistical measures from finite data grows rapidly with dimensionality, making the dimensionality of the problem a priority for any networks containing more than a couple nodes.

One approach for network inference is to test each candidate causal link conditioned on all other variables [32]. That is, a direct link  $j \rightarrow i$  is inferred if such a relationship remains effective when conditioning on all other variables in the system. Although intuitive and correct in theory, this method requires computing a statistic in a sample space as high dimensional as the entire system and therefore falls short when applied to a large networks. The PC algorithm [66] overcomes this difficulty by repeated testing of the candidate causal link conditioned on subsets of the remaining variables [60, 61]. To be more specific, a link  $j \rightarrow i$  is disqualified as a candidate causal relationship if it is insignificant when conditioned on some subset of the nodes. The advantage of the PC algorithm is that it reduces the dimensionality of the sample space the test of independence to be proportional to the size of the conditioning set (which in some cases can be much smaller than the system size). However, unless the maximum degree of the nodes are known *a priori*, the algorithm in principle needs to be performed for combinations of subsets as the conditioning sets up to the size of the entire network. In this respect, regardless of the dimensionality of the sample space, the combinatorial search itself can be computationally infeasible for moderate to large networks. In practice, tradeoff needs to be made between an algorithm’s computational cost and data efficiency (in terms of the estimation of the test statistic).

In this paper we develop theory of causation entropy—a type of conditional mutual information designed for causal network inference. In particular, we prove the optimal causation entropy principle for Markov processes: the set of nodes that directly cause a given node is the unique minimal set of nodes that maximizes causation entropy. This principle allows us to convert the problem of causality inference into the

optimization of causation entropy. We further show that this optimization problem, which appears to be combinatorial, can in fact be solved by simple greedy algorithms, which are both computationally efficient and data efficient. We verify the effectiveness of the proposed algorithms through analytical and numerical investigations of Gaussian processes on various network types including trees, loops, and random networks. Somewhat surprisingly, our results suggest that it is the density of links and information diffusion rate rather than the number of nodes in a network that determines the minimal sample size required for accurate inference.

**2. Stochastic Process and Causal Network Inference.** We begin by introducing a theoretical framework for inferring causal networks from high-dimensional time series. This framework is general in that it is applicable to both linear and non-linear systems with or without added noise.

Consider a network (graph)  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , with  $\mathcal{V} = \{1, 2, \dots, n\}$  being the set of nodes and  $\mathcal{E} \subset \mathcal{V} \times \mathcal{V} \times \mathbb{R}$  being the set of weighted links (or edges). The *adjacency matrix*  $A = [A_{ij}]_{n \times n}$  is defined as

$$(2.1) \quad A_{ij} = \begin{cases} \text{weight of the link } j \rightarrow i, & \text{if } j \rightarrow i \text{ in the network;} \\ 0, & \text{otherwise.} \end{cases}$$

We use  $\chi_0(A)$  to denote the corresponding unweighted adjacency matrix defined entry-wise by  $\chi_0(A)_{ij} = 1$  iff  $A_{ij} \neq 0$  and  $\chi_0(A)_{ij} = 0$  iff  $A_{ij} = 0$ . We define the set of *causal parents* of  $i$  as

$$(2.2) \quad N_i = \{j | A_{ij} \neq 0\} = \{j | \chi_0(A)_{ij} = 1\}.$$

For a subset of nodes  $I \subset \mathcal{V}$ , we similarly define its set of causal parents as

$$(2.3) \quad N_I = \cup_{i \in I} N_i.$$

We consider stochastic network dynamics in the following form (for each node  $i$ )

$$(2.4) \quad X_t^{(i)} = f_i(A_{i1}X_{t-1}^{(1)}, A_{i2}X_{t-1}^{(2)}, \dots, A_{ij}X_{t-1}^{(j)}, \dots, A_{in}X_{t-1}^{(n)}, \xi_t^{(i)})$$

where  $X_t^{(i)} \in \mathbb{R}^d$  is a random variable representing the state of node  $i$  at time  $t$ ,  $\xi_t^{(i)} \in \mathbb{R}^d$  is the random fluctuation on node  $i$  at time  $t$ , and  $f_i : \mathbb{R}^{d \times (n+1)} \rightarrow \mathbb{R}^d$  models the functional dependence of the state of node  $i$  on the past states of nodes  $j$  with  $A_{ij} \neq 0$ . Note that other than the noise term  $\xi_t^{(i)}$ , the state  $X_t^{(i)}$  only depends (stochastically) on the past states of its causal parents,  $X_{t-1}^{(j)}$  ( $j \in N_i$ ).

For a subset  $K = \{k_1, k_2, \dots, k_q\} \subset \mathcal{V}$ , we define

$$(2.5) \quad X_t^{(K)} \equiv [X_t^{(k_1)}, X_t^{(k_2)}, \dots, X_t^{(k_q)}]^\top.$$

If  $K = \mathcal{V}$ , we simplify the notation and denote

$$(2.6) \quad X_t \equiv X_t^{(\mathcal{V})} = [X_t^{(1)}, X_t^{(2)}, \dots, X_t^{(n)}]^\top.$$

**2.1. Problem of Causal Network Inference and Challenges.** Given quantitative observations of the dynamic states of individual nodes, often in the form of time series, a central problem is to infer its (causal) system dynamics, which involves the inference of (1) the causal network topology,  $\chi_0(A)$ ; (2) the link weights,  $\{A_{ij}\}$ ;

and (3) the specific forms of functional dependence between nodes,  $\{f_i\}$ . These problems are interrelated and all challenging. We focus on the first problem: inferring the causal network topology  $\chi_0(A)$ , which serves as the skeleton of the actual network dynamics. See Fig. 1.1 as a schematic illustration. In particular, the problem of causal network inference can be casted mathematically as:

$$(2.7) \quad \begin{cases} \text{Given:} & \text{Samples of the node states } x_t^{(i)} \text{ (} i = 1, 2, \dots, n; t = 1, 2, \dots, T\text{).} \\ \text{Goal:} & \text{Infer the structure of the underlying causal network,} \\ & \text{i.e., find } \operatorname{argmin}_{\hat{A}} \|\chi_0(A) - \hat{A}\|_0, \text{ where } \|M\|_0 \equiv \sum_{i,j} |M_{ij}|^0. \end{cases}$$

One key challenge is that in many applications, the number of nodes  $n$  is often *large* (usually hundreds at least), but the sample size  $T$  is much smaller than needed for reliable estimation of the  $(n \times d)$ -dimensional joint distribution. We propose that a practical causation inference method should fulfill the following three requirements:

1. *Model-free.* The method should not rely on assumptions about either the form or parameters of a model that underlie the process.
2. *Computational Efficient.* The method should be computationally efficient.
3. *Data Efficient.* The method should achieve high accuracy with relatively small number of samples (i.e., convergence in probability needs to be fast).

In this paper we address the model-free requirement by utilizing information-theoretic measures, and in particular, by using causation entropy. On the other hand, our theoretical developments of the optimal causation entropy principle enables us to develop algorithms that are both computationally efficient and data efficient.

**2.2. Markov Assumptions.** We study the system in a probabilistic framework assuming stationarity and existence of a continuous distribution. We further make the following assumptions regarding the conditional distributions  $p(\cdot|\cdot)$  arising from the stationary process given by Eq. (2.4). For every node  $i \in \mathcal{V}$  and time indices  $t, t'$ :

$$(2.8) \quad \begin{cases} (1) \text{ Temporally Markov:} \\ \quad p(X_t|X_{t-1}, X_{t-2}, \dots) = p(X_t|X_{t-1}) = p(X_{t'}|X_{t'-1}). \\ (2) \text{ Spatially Markov:} \\ \quad p(X_t^{(i)}|X_{t-1}) = p(X_t^{(i)}|X_{t-1}^{(N_i)}). \\ (3) \text{ Faithfully Markov:} \\ \quad p(X_t^{(i)}|X_{t-1}^{(K)}) \neq p(X_t^{(i)}|X_{t-1}^{(L)}) \text{ whenever } (K \cap N_i) \neq (L \cap N_i). \end{cases}$$

Throughout the paper, the relationship between two probability density functions  $p_1$  and  $p_2$  are denoted as “ $p_1 = p_2$ ” iff they equal *almost everywhere*, and “ $p_1 \neq p_2$ ” iff there is a set of positive measure on which the two functions do not equal.

In Eq. (2.8), condition (1) states that the underlying dynamics is a time-invariant Markov process<sup>1</sup>. Condition (2) is often referred to as the (local) Markov property [44], which we call Spatially Markov here to differ from Temporally Markov. This condition guarantees that in determining the future state of a node, if knowledge about the past states of all its causal parents  $N_i$  [as defined in Eq. (2.2)] is given, information about the past of any other node becomes irrelevant. Finally, condition (3) ensures that the

<sup>1</sup>If the process is Markov but with higher order, our approach is to convert it into a first-order one as illustrated in Appendix A and then apply the theory and algorithms in the main body of the paper to the resulting first-order process.

set of causal parents is unique and that every causal parent presents an observable effect regardless of the information about other causal parents<sup>2</sup>.

The conditional independence between two random variables  $X$  and  $Y$  given  $Z$  is denoted by  $(X \perp\!\!\!\perp Y \mid Z)$ , i.e.,

$$(2.9) \quad (X \perp\!\!\!\perp Y \mid Z) \iff p(X|Y, Z) = p(X|Z).$$

The following results regarding conditional independence will be useful in later sections and are direct consequences of the basic axioms of probability theory [31, 44, 53]:

$$(2.10) \quad \left\{ \begin{array}{l} \text{Symmetry: } (X \perp\!\!\!\perp Y \mid Z) \iff (Y \perp\!\!\!\perp X \mid Z). \\ \text{Decomposition: } (X \perp\!\!\!\perp YW \mid Z) \implies (X \perp\!\!\!\perp Y \mid Z). \\ \text{Weak union: } (X \perp\!\!\!\perp YW \mid Z) \implies (X \perp\!\!\!\perp Y \mid ZW). \\ \text{Contraction: } (X \perp\!\!\!\perp Y \mid Z) \wedge (X \perp\!\!\!\perp W \mid ZY) \implies (X \perp\!\!\!\perp YW \mid Z). \\ \text{Intersection: } (X \perp\!\!\!\perp Y \mid ZW) \wedge (X \perp\!\!\!\perp W \mid ZY) \implies (X \perp\!\!\!\perp YW \mid Z). \end{array} \right.$$

Here “ $\wedge$ ” denotes the logical operations “and” (the symbol “ $\vee$ ” is used later for “or”), and  $YW$  denotes a joint random variable of  $Y$  and  $W$ .

**2.3. Causation Entropy as an Information-Theoretic Measure.** We review several fundamental concepts in information theory, leading to causation entropy, a model-free information-theoretic statistic that can be used to infer direct causal relationships [71].

Originally proposed by Shannon as a measure of uncertainty and complexity, the (differential) *entropy* of a continuous random variable  $X \in \mathbb{R}^n$  is defined as [16, 64]<sup>3</sup>

$$(2.11) \quad h(X) = - \int p(x) \log p(x) dx,$$

where  $p(x)$  is the probability density function of  $X$ . The joint and conditional entropies between two random variables  $X$  and  $Y$  are defined as [also see Fig. 2.1(a)]

$$(2.12) \quad \left\{ \begin{array}{l} \text{Joint entropy: } h(X, Y) \equiv h(Y, X) \equiv - \int p(x, y) \log p(x, y) dx dy. \\ \text{Conditional entropies: } \begin{cases} h(X|Y) \equiv - \int p(x, y) \log p(x|y) dx dy; \\ h(Y|X) \equiv - \int p(x, y) \log p(y|x) dx dy. \end{cases} \end{array} \right.$$

For more than two random variables, the entropies are similarly defined (as above) by grouping the variables into two classes, one acting as  $X$  and the other as  $Y$ .

The *mutual information* between two random variables  $X$  and  $Y$  (conditioning on  $Z$ ) can be interpreted as a measure of the deviation from independence between  $X$  and  $Y$  (conditioning on  $Z$ ). The corresponding unconditioned and conditional mutual information are defined respectively as

$$(2.13) \quad \left\{ \begin{array}{l} \text{Mutual information: } I(X; Y) \equiv h(X) - h(X|Y) \equiv h(Y) - h(Y|X). \\ \text{Conditional mutual information:} \\ \quad I(X; Y|Z) \equiv h(X|Z) - h(X|Y, Z) \equiv h(Y|Z) - h(Y|X, Z). \end{array} \right.$$

<sup>2</sup>Note that without condition (3), the “true positive” statement in Theorem 2.2 is no longer valid. One simple example is given in Appendix B to illustrate this point.

<sup>3</sup>We follow the convention in Ref. [16] to use  $h(\cdot)$  for the entropy of a continuous random variable and reserve  $H(\cdot)$  for the entropy of a discrete random variable. In the discrete case, we need to replace the integral by summation and probability density by probability mass function in the definition.

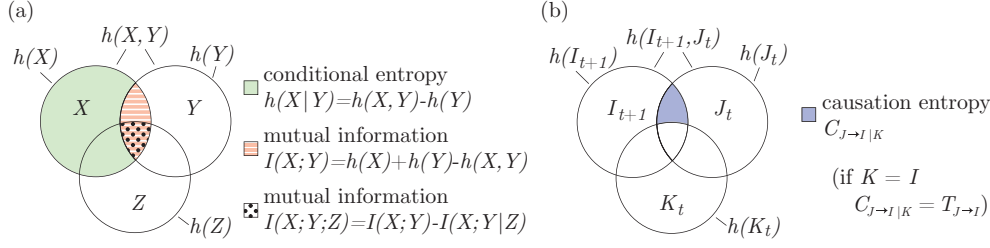


FIG. 2.1. Venn-like diagrams for information-theoretic measures. (a) Visualization of the relationships between entropy, conditional entropy, and mutual information. (b) Visualization of the relationships between conditional entropy, causation entropy, and transfer entropy. In the picture of (b), letters  $I$ ,  $J$ , and  $K$  are used to denote  $X^{(I)}$ ,  $X^{(J)}$ , and  $X^{(K)}$ , respectively.

The mutual information among three variables  $X$ ,  $Y$ , and  $Z$  is defined as<sup>4</sup>

$$(2.14) \quad I(X; Y; Z) \equiv I(X; Y) - I(X; Y|Z) \equiv I(Y; Z) - I(Y; Z|X) \equiv I(X; Z) - I(X; Z|Y),$$

The mutual information between two variables is always nonnegative,  $I(X; Y) \geq 0$ , with equality if and only if  $X$  and  $Y$  are independent. Similarly,  $I(X; Y|Z) \geq 0$ , with equality if and only if  $X$  and  $Y$  are independent when conditioned on  $Z$ . Interestingly, for three or more variables, such an inequality does not hold: the mutual information  $I(X; Y; Z)$  can be either positive, negative or zero [47]. Figure 2.1(a) visualizes the relationships between entropy, conditional entropy, and mutual information.

To measure the directionality of information flow between two random processes, Schreiber proposed a specific type of conditional mutual information called *transfer entropy* [63]. For a stationary first-order Markov process such as the one given by Eq. (2.4), the transfer entropy from  $j$  to  $i$  can be expressed as

$$(2.15) \quad T_{j \rightarrow i} \equiv h(X_{t+1}^{(i)} | X_t^{(i)}) - h(X_{t+1}^{(i)} | X_t^{(i)}, X_t^{(j)}),$$

where  $h(\cdot|\cdot)$  denotes conditional entropy [16]. Since  $h(X_{t+1}^{(i)} | X_t^{(i)})$  measures the uncertainty of  $X_{t+1}^{(i)}$  given information about  $X_t^{(i)}$  and  $h(X_{t+1}^{(i)} | X_t^{(i)}, X_t^{(j)})$  measures the uncertainty of  $X_{t+1}^{(i)}$  given information about *both*  $X_t^{(i)}$  and  $X_t^{(j)}$ , the transfer entropy  $T_{j \rightarrow i}$  can be interpreted as the *reduction of uncertainty* about future states of  $X^{(i)}$  when the current state of  $X^{(j)}$  is provided *in addition to that of*  $X^{(i)}$ .

Networks of practical interest inevitably contain (many) more than two nodes. As we will show later, without appropriate conditioning transfer entropy fails to distinguish between direct and indirect causality in networks. To overcome the pairwise limitation of transfer entropy, we define *causation entropy*. The relationships between entropy, transfer entropy and causation entropy are illustrated in Fig 2.1(b).

DEFINITION 2.1 (Causation Entropy [71]). *The causation entropy from the set*

<sup>4</sup>This quantity is often referred to as *interaction information* [47] or *co-information* [7]. Another multivariate generalizations of mutual information is total correlation [78] (also known as multivariate constraint [24] or multi-information [69]).

of nodes  $J$  to the set of nodes  $I$  conditioning on the set of nodes  $K$  is defined as<sup>5</sup>

$$(2.16) \quad C_{J \rightarrow I|K} = h(X_{t+1}^{(I)}|X_t^{(K)}) - h(X_{t+1}^{(I)}|X_t^{(K)}, X_t^{(J)}),$$

where  $I, J, K$  are all subset of  $\mathcal{V} = \{1, 2, \dots, n\}$ . In particular, if  $J = \{j\}$  and  $I = \{i\}$ , we simplify the notation as  $C_{j \rightarrow i|K}$ . If the conditioning set  $K = \emptyset$ , we often omit it and simply write  $C_{J \rightarrow I}$ .

REMARK 2.1. Causation entropy is a natural generalization of transfer entropy from measuring pairwise causal relationships to network relationships of many variables. In particular, if  $j \in K$ , then the causation entropy  $C_{j \rightarrow i|K} = 0$  as  $j$  does not carry extra information (compared to that of  $K$ ). On the other hand, if  $K = \{i\}$ , causation entropy recovers transfer entropy, i.e.,

$$(2.17) \quad C_{j \rightarrow i|i} = T_{j \rightarrow i}.$$

Interestingly, in this framework we see that transfer entropy assumes that nodes are self-causal, whereas causation entropy relaxes this assumption. Preliminary exploration of the differences between the two measures can be found in Ref. [71].

REMARK 2.2. We note that in addition to Ref. [71], the conditional mutual information between time-lagged variables has been proposed as a statistic for network inference in a few previous studies [21, 60, 61, 75] (although not referred to as transfer or causation entropy).

REMARK 2.3. It seems plausible to conjecture that if two subsets of the nodes satisfy  $K_1 \subset K_2$ , then  $C_{j \rightarrow i|K_1}$  would be no less than  $C_{j \rightarrow i|K_2}$ . We remark that this statement about monotonicity is false (see the two examples below).

Example 1. Consider the stochastic process

$$(2.18) \quad X_t^{(1)} = X_{t-1}^{(2)} + X_{t-1}^{(3)}$$

where  $X_t^{(k)}$  are i.i.d Bernoulli variables:  $P(X_t^{(k)} = 0) = P(X_t^{(k)} = 1) = 0.5$  ( $k = 2, 3$ ). Let  $i = 1, j = 2, K_1 = \emptyset$  and  $K_2 = \{3\}$ . It follows that

$$(2.19) \quad \begin{cases} C_{2 \rightarrow 1|\emptyset} = \frac{3}{2} \log 2 - \log 2 = \frac{1}{2} \log 2 \\ C_{2 \rightarrow 1|\{3\}} = \log 2 - 0 = \log 2 \end{cases} \Rightarrow C_{2 \rightarrow 1|\emptyset} < C_{2 \rightarrow 1|\{3\}}.$$

Example 2. Consider the stochastic process

$$(2.20) \quad X_{t+1}^{(1)} = X_t^{(3)}, X_{t+1}^{(2)} = X_t^{(3)},$$

where  $X_t^{(3)}$  are Bernoulli variables with  $P(X_t^{(3)} = 0) = P(X_t^{(3)} = 1) = 0.5$ . Let  $i = 1, j = 2, K_1 = \emptyset$  and  $K_2 = \{3\}$ . It follows that

$$(2.21) \quad \begin{cases} C_{2 \rightarrow 1|\emptyset} = \log 2 - 0 = \log 2 \\ C_{2 \rightarrow 1|\{3\}} = 0 - 0 = 0 \end{cases} \Rightarrow C_{2 \rightarrow 1|\emptyset} > C_{2 \rightarrow 1|\{3\}}.$$

The seemingly paradoxical observation that  $C_{j \rightarrow i|K_1}$  can either be larger or smaller than  $C_{j \rightarrow i|K_2}$  despite the fact that  $K_1 \subset K_2$  can be understood as follows: When

<sup>5</sup> Note that the definitions in Eq. (2.15) and Eq. (2.16) can be extended for asymptotically stationary processes by taking the limit of  $t \rightarrow \infty$ , although the proofs in this paper do not directly apply to such general scenario.



$K_1 \subset K_2$ ,  $C_{j \rightarrow i|K_1} - C_{j \rightarrow i|K_2}$  corresponds to the mutual information among the three variables  $X_{t+1}^{(i)}|X_t^{(K_1)}$ ,  $X_{t+1}^{(i)}|X_t^{(j)}$  and  $X_{t+1}^{(i)}|X_t^{(K_2-K_1)}$  (see Fig. 2.1). Contrary to the two-variable case where mutual information is always nonnegative, the mutual information among three (or more) variables can either be positive, negative or zero [47].

**2.4. Theoretical Properties of Causation Entropy and the Optimal Causation Entropy Principle.** In the following we show that analysis of causation entropy leads to exact network inference for the network stochastic process given by Eq. (2.4) subject to the Markov assumptions in Eq. (2.8).

We start by exploring basic analytical properties of causation entropy, which is presented as Theorem 2.2 and also summarized in Fig. 2.2.

**THEOREM 2.2** (Basic analytical properties of causation entropy). *Suppose that the network stochastic process given by Eq. (2.4) satisfies the Markov assumptions in Eq. (2.8). Let  $I \subset \mathcal{V}$  be a set of nodes and  $N_I$  be its causal parents. Consider two sets of nodes  $J \subset \mathcal{V}$  and  $K \subset \mathcal{V}$ . The following results hold:*

- (a) (Redundancy) If  $J \subset K$ , then  $C_{J \rightarrow I|K} = 0$ .
- (b) (No false positive) If  $N_I \subset K$ , then  $C_{J \rightarrow I|K} = 0$  for any set of nodes  $J$ .
- (c) (True positive) If  $J \subset N_I$  and  $J \not\subset K$ , then  $C_{J \rightarrow I|K} > 0$ .
- (d) (Decomposition)  $C_{J \rightarrow I|K} = C_{(K \cup J) \rightarrow I} - C_{K \rightarrow I}$ .

*Proof.* Under the Temporal Markov Condition in Eq. (2.8), there is no time dependence of the distributions. For notational simplicity we denote the joint distribution  $p(X_{t+1}^{(I)} = i, X_t^{(J)} = j, X_t^{(K)} = k)$  by  $p(i, j, k)$  and use similar notation for the marginal and conditional distributions. It follows that

$$\begin{aligned} C_{J \rightarrow I|K} &= h(X_{t+1}^{(I)}|X_t^{(K)}) - h(X_{t+1}^{(I)}|X_t^{(K)}, X_t^{(J)}) = - \int p(i, j, k) \log \left[ \frac{p(i|k)}{p(i|j, k)} \right] didjdk \\ &\geq - \log \int p(i, j, k) \frac{p(i|k)}{p(i|j, k)} didjdk \quad (\text{by Jensen's inequality [58]}) \\ (2.22) \quad &= - \log \int p(j, k) \frac{p(i, k)}{p(k)} didjdk = - \log(1) = 0, \end{aligned}$$

where equality holds if and only if  $p(i|k) = p(i|j, k)$  almost everywhere. The above inequality is also known as the *Gibbs' inequality* in statistical physics [26].

To prove (a), we note that  $J \subset K$  implies that  $p(i|k) = p(i|j, k)$  and therefore equality holds (rather than inequality) in Eq. (2.22).

To prove (b), it suffices to show that for  $J \not\subset K$ ,  $C_{J \rightarrow I|K} = 0$ . Since  $J \not\subset K$  and  $N_I \subset K$ , based on the Spatial Markov Condition in Eq. (2.8), we have:

$$(2.23) \quad p(X_{t+1}^{(I)}|X_t) = p(X_{t+1}^{(I)}|X_t^{(K \cup J)}) = p(X_{t+1}^{(I)}|X_t^{(K)}) = p(X_{t+1}^{(I)}|X_t^{(N_I)}).$$

Therefore  $p(i|j, k) = p(i|k)$  and equality holds in Eq. (2.22).

To prove (c), we use the Faithfully Markov Condition in Eq. (2.8). Since  $J \subset N_I$  and  $J \not\subset K$ , it follows that

$$(2.24) \quad p(X_{t+1}^{(I)}|X_t^{(K)}) = p(X_{t+1}^{(I)}|X_t^{(K \cap N_I)}) \neq p(X_{t+1}^{(I)}|X_t^{(K)}, X_t^{(J)}).$$

Thus,  $p(i|j, k) \neq p(i|k)$  and strictly inequality holds in Eq. (2.22).

Finally, part (d) follows directly from the definition of  $C$ .  $\square$

Theorem 2.2 allows us to convert the problem of causal network inference into the problem of estimating causation entropy among nodes. In particular, for a given

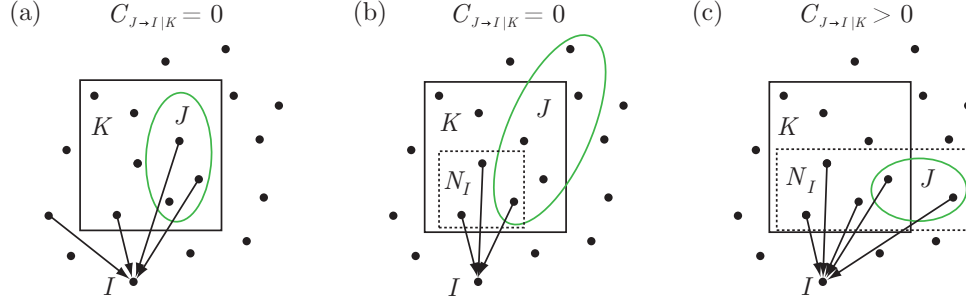


FIG. 2.2. Basic analytical properties of causation entropy (Theorem 2.2) allowing for the inference of the causal parents  $N_I$  of a set of nodes  $I$ . (a) Redundancy: If  $J$  is a subset of the conditioning set  $K$  ( $J \subset K$ ), then the causation entropy  $C_{J \rightarrow I|K} = 0$ . (b) No false positive: If  $N_I$  is already included in the conditioning set  $K$  ( $N_I \subset K$ ), then  $C_{J \rightarrow I|K} = 0$ . (c) True positive: If a set  $J$  contains at least one causal parent of  $I$  that does not belong to the conditioning set  $K$ , i.e.,  $(J \subset N_I) \wedge (J \not\subset K)$ , then  $C_{J \rightarrow I|K} > 0$ .

set of nodes  $I$ , each node  $j$  can *in principle* be checked independently to determine whether or not it is a causal parent of  $I$  via either of the following two equivalent criteria (proved in Theorem 2.3(a) below)

$$(2.25) \quad \begin{cases} (1) \text{ Node } j \in N_I \text{ iff there is a set } K \supset N_I, \text{ such that } C_{j \rightarrow I|(K-\{j\})} > 0; \\ (2) \text{ Node } j \in N_I \text{ iff for any set } K \subset \mathcal{V}, C_{j \rightarrow I|(K-\{j\})} > 0. \end{cases}$$

Practical application of either criteria to infer large networks is challenging. Criterion (1) requires a conditioning set  $K$  that contains  $N_I$  as its subset. Since  $N_I$  is generally unknown, one often must use  $K = \mathcal{V}$ . When the network is large ( $n \gg 1$ ), this requires the estimation of causation entropy for very high dimensional random variables from limited data, which is inherently unreliable [60, 61]. Criterion (2), on the other hand, requires a combinatorial search over all subsets making it computationally infeasible.

In the following we prove the two inference criteria in Eq. (2.25). Furthermore, we show that *the set of causal parents is the minimal set of nodes that maximizes causation entropy*, which we refer to as the *optimal causation entropy principle*.

**THEOREM 2.3** (Optimal causation entropy principle for causal network inference). *Suppose that the network stochastic process given by Eq. (2.4) satisfies the Markov properties in Eq. (2.8). Let  $I \subset \mathcal{V}$  be a given set of nodes and  $N_I$  be the set of  $I$ 's causal parents, as defined in Eq. (2.3). It follows that*

- (a) (Direct inference) Node  $j \in N_I$  iff  $\Leftrightarrow \exists K \supset N_I$  such that  $C_{j \rightarrow I|(K-\{j\})} > 0 \Leftrightarrow \forall K \subset \mathcal{V}, C_{j \rightarrow I|(K-\{j\})} > 0$ .
- (b) (Partial conditioning removal) If there exists  $K \subset \mathcal{V}$  such that  $C_{j \rightarrow I|(K-\{j\})} = 0$ , then  $j \notin N_I$ .
- (c) (Optimal causation entropy principle) The set of causal parents is the minimal set of nodes with maximal causation entropy.

Define the family of sets with maximal causation entropy as

$$(2.26) \quad \mathcal{K} = \{K | \forall K' \subset \mathcal{V}, C_{K' \rightarrow I} \leq C_{K \rightarrow I}\}.$$

Then the set of causal parents satisfies

$$(2.27) \quad N_I = \cap_{K \in \mathcal{K}} K = \operatorname{argmin}_{K \in \mathcal{K}} K.$$

*Proof.* First we prove part (a). If  $j \in N_I$ , then for every  $K \subset \mathcal{V}$ ,  $C_{j \rightarrow I|(K-\{j\})} > 0$  following Theorem 2.2(c). This proves both “ $\Rightarrow$ ”. On the other hand, suppose that  $\forall K \subset \mathcal{V}, C_{j \rightarrow I|(K-\{j\})} > 0$ , then for  $K = \mathcal{V} \supset N_I$ , it follows that  $C_{j \rightarrow i|(\mathcal{V}-\{j\})} > 0$ . Node  $j \in N_I$  since otherwise  $(\mathcal{V}-\{j\}) \supset N_I$  which would imply that  $C_{j \rightarrow i|(\mathcal{V}-\{j\})} = 0$  from Theorem 2.2(b). Therefore, the two “ $\Leftarrow$ ”s are also proven.

Next, part (b) follows directly from the contrapositive of Theorem 2.2(c).

Finally, we prove part (c). Note that if  $N_I \not\subset K$ , then  $J = N_I - K \neq \emptyset$ , and so  $C_{(K \cup J) \rightarrow I} - C_{K \rightarrow I} = C_{J \rightarrow I|K} > 0$ . Therefore,  $K \in \mathcal{K} \Rightarrow N_I \subset K$ . This implies  $N_I \subset \cap_{K \in \mathcal{K}} K$ . On the other hand, if  $\exists j \in \cap_{K \in \mathcal{K}} K$  with  $j \notin N_I$ . Let  $K \in \mathcal{K}$  and  $L = K - \{j\}$ . Since  $j \notin N_I$ , we have  $N_I \subset L \subset K$ , and therefore  $C_{K \rightarrow I} - C_{L \rightarrow I} = C_{j \rightarrow I|L} = 0$ , where the second equality follows from Theorem 2.2(c). This shows that  $L \in \mathcal{K}$ , contradicting with  $j \in \cap_{K \in \mathcal{K}} K$ . So  $j \in \cap_{K \in \mathcal{K}} K \Rightarrow j \in N_I$ , which implies that  $\cap_{K \in \mathcal{K}} K \subset N_I$ . Since  $\mathcal{K}$  is finite, it follows that  $\cap_{K \in \mathcal{K}} K = \operatorname{argmin}_{K \in \mathcal{K}} K$ .  $\square$

Based on the optimal causation entropy principle, it seems straightforward to solve the minimax optimization for the inference of  $N_I$  by enumerating all subsets of  $\mathcal{V}$  with increasing cardinality (starting from  $\emptyset$ ), and terminating when a set  $K$  is found to be have maximal causation entropy among all subsets of cardinality  $|K| + 1$  (i.e., adding any node  $j$  to set  $K$  does not increase the causation entropy  $C_{K \rightarrow I}$ ). Based on Theorem 2.3, the set  $K = N_I$ . However, this brute-force approach requires  $\mathcal{O}(n^{|N_I|})$  causation entropy evaluations, which is computationally inefficient and therefore infeasible for the inference of real world networks which often contain large number of nodes ( $n \gg 1$ ). Such limitation is removed only when the number of causal parents is moderately small,  $|N_I| = \mathcal{O}(1)$ . In the following section we develop additional theory and algorithms to efficiently solve this minimax optimization problem for causal network inference.

**2.5. Computational Causal Network Inference.** Algorithmically, causal network inference via the optimal causation entropy principle should require as few computations as necessary (computational efficiency) and as few data samples as possible while retaining accuracy (data efficiency). We introduce two such algorithms that jointly infer the causal network. For a given node  $i$ , the goal is to infer its causal parents, as illustrated by nodes in the shaded region of Fig. 2.3(a). Algorithm 2.1 aggregatively identifies nodes that form a superset of the causal parents,  $K \supset N_i$  (proven by Lemma 2.4, illustrated in Fig. 2.3(b)). Start from a set  $K \supset N_i$ , Algorithm 2.2 prunes away non-causal nodes from  $K$  leaving only the causal parents  $N_i$  (proven by Lemma 2.5, illustrated in Fig. 2.3(c)).

LEMMA 2.4 (Aggregative Discovery of Causal Nodes). *Suppose that the network stochastic process given by Eq. (2.4) satisfies the Markov properties in Eq. (2.8). Let  $I \subset \mathcal{V}$  and  $N_I$  be its causal parents. Define the sequences of numbers  $\{x_1, x_2, \dots\}$ , nodes  $\{p_1, p_2, \dots\}$ , and nested sets  $\{K_0, K_1, K_2, \dots\}$  as:  $K_0 = \emptyset$ , and*

$$(2.28) \quad \begin{cases} x_i = \max_{x \in (\mathcal{V}-K_{i-1})} C_{x \rightarrow I|K_{i-1}}, \\ p_i = \operatorname{argmax}_{x \in (\mathcal{V}-K_{i-1})} C_{x \rightarrow I|K_{i-1}}, \\ K_i = \{p_1, p_2, \dots, p_i\} \end{cases}$$

for every  $i \geq 1$ . There exists a number  $q$ , with  $|N_I| \leq q \leq n$ , such that

- (a) The numbers  $x_i > 0$  for  $1 \leq i \leq q$  and  $x_i = 0$  for  $i > q$ .
- (b) The set of causal parents  $N_I \subset K_q = \{x_1, x_2, \dots, x_q\}$ .

*Proof.* If  $N_I = \emptyset$ , the lemma holds trivially. Suppose that  $|N_I| \geq 1$  and so  $x_1 > 0$ .

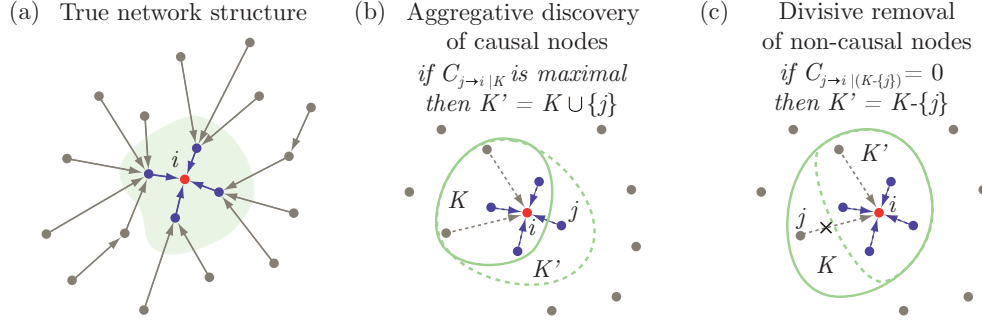


FIG. 2.3. Causal network inference by optimal causation entropy. (a) Causal parents and non-causal nodes of a node  $i$ . Causal network inference corresponds to identifying the causal parents  $N_i$  (nodes in shaded region) for every node  $i \in \mathcal{V}$ . (b) Nodes are added to the set  $K$  in an aggregative fashion, maximizing causation entropy at each step (see Algorithm 2.1). (c) Starting from a set  $K \supset N_i$  ( $K$  obtained by Algorithm 2.1), non-causal nodes are progressively removed from  $K$  if their causation entropy to node  $i$  conditioned on the rest of  $K$  is zero (see Algorithm 2.2).

---

### Algorithm 2.1 Aggregative Discovery of Causal Nodes

---

**Input:** Set of nodes  $I \subset \mathcal{V}$

**Output:**  $K$  (which will include  $N_I$  as its subset)

- 1: Initialize:  $K \leftarrow \emptyset$ ,  $x \leftarrow \infty$ ,  $p \leftarrow \emptyset$ .
  - 2: **while**  $x > 0$  **do**
  - 3:    $K \leftarrow K \cup \{p\}$
  - 4:   **for** every  $j \in (\mathcal{V} - K)$  **do**
  - 5:      $x_j \leftarrow C_{j \rightarrow I|K}$
  - 6:   **end for**
  - 7:    $x \leftarrow \max_{j \in (\mathcal{V} - K)} x_j$ ,  $p \leftarrow \operatorname{argmax}_{j \in (\mathcal{V} - K)} x_j$
  - 8: **end while**
- 

To prove (a), we define  $q \equiv \min_{x_i=0}(i-1)$  (if all  $x_i > 0$ , define  $q \equiv n$ ). By construction,  $x_i > 0$  when  $i \leq q$  and  $x_{q+1} = 0$ . This implies that  $N_I \subset K_q$  since otherwise there is a node  $j$  with  $C_{j \rightarrow I|K_q} > 0 \Rightarrow x_{q+1} > 0$ . For any  $i > q$ ,  $N_I \subset K_q \subset K_{i-1}$ , and thus  $C_{j \rightarrow I|K_{i-1}} = 0$  for all  $j \in (\mathcal{V} - K_{i-1})$ , which implies that  $x_i = 0$ .

To prove (b), we note that if there is a node  $j \in N_I$  such that  $j \notin K_q$ , then by the definition of  $x_i$  and Theorem 2.2(c), it follows that  $x_{q+1} \geq C_{j \rightarrow I|K_q} > 0$ . This is in contradiction with the fact that  $x_i = 0$  for all  $i > q$ . Therefore,  $N_I \subset K_q$ .  $\square$

Algorithm 2.1 recursively constructs the set  $K_q \supset N_I$  (further denoted as  $K$ ) as described by Lemma 2.4 and illustrated in Fig. 2.3(b). To remove indirect and spurious nodes in  $K$  that do not belong to  $N_I$ , we apply the result of Theorem 2.2(c),  $C_{j \rightarrow I|(K-\{j\})} = 0 \Rightarrow j \notin N_I$ . This gives rise to Lemma 2.5 and Algorithm 2.2.

LEMMA 2.5 (Progressive Removal of Non-Causal Nodes). *Suppose that the network stochastic process given by Eq. (2.4) satisfies the Markov properties in Eq. (2.8). Let  $I \subset \mathcal{V}$  and  $N_I$  be its causal parents. Let  $K = \{p_1, p_2, \dots, p_q\}$  such that  $K \supset N_I$ . Define the sequence of sets  $\{K_0, K_1, K_2, \dots, K_q\}$  by  $K_0 = K$ , and*

$$(2.29) \quad K_i = \begin{cases} K_{i-1}, & \text{if } C_{p_i \rightarrow I|(K_{i-1}-\{p_i\})} > 0; \\ K_{i-1} - \{p_i\}, & \text{if } C_{p_i \rightarrow I|(K_{i-1}-\{p_i\})} = 0. \end{cases}$$

**Algorithm 2.2 Progressive Removal of Non-Causal Nodes****Input:** Sets of nodes  $I \subset \mathcal{V}$  and  $K \subset \mathcal{V}$ **Output:**  $\hat{N}_I$  (inferred set of causal parents of  $I$ )

---

```

1: for every  $j \in K$  do
2:   if  $C_{j \rightarrow I|(K-\{j\})} = 0$  then
3:      $K \leftarrow K - \{j\}$ 
4:   end if
5: end for
6:  $\hat{N}_I \leftarrow K$ 

```

---

for every  $1 \leq i \leq q$ . Then  $K_q = N_I$ .

*Proof.* By definition,  $K_0 = K \supset N_I$ . We prove that  $K_q \supset N_I$  by induction. Suppose that  $K_{i-1} \supset N_I$ . If node  $p_i \in N_I$ , then  $C_{p_i \rightarrow I|(K_{i-1}-\{p_i\})} > 0$  by Theorem 2.2(c) and therefore  $K_i = K_{i-1} \supset N_I$ . If node  $p_i \notin N_I$ , then  $K_i \supset K_{i-1} - \{p_i\} \supset N_I$ .

Next we prove that  $K_q \subset N_I$ . Suppose that node  $p_i \notin N_I$ . Since  $K_{i-1} \supset N_I$ , the causation entropy  $C_{p_i \rightarrow I|(K_{i-1}-\{p_i\})} = 0$  by Theorem 2.2(b), and so  $K_i = K_{i-1} - \{p_i\}$ . Therefore,  $p \notin K_i \supset K_q$ , which implies that  $K_q \subset N_I$  (contrapositive).  $\square$

Algorithm 2.2 iteratively removes nodes that are not causal parents from a set  $K$  until the set converges to  $N_I$  as described by Lemma 2.5 and illustrated in Fig. 2.3(c).

Jointly, Algorithms 2.1 and 2.2 can be applied to identify the causal parents of each node, thus inferring the entire causal network<sup>6</sup>.

**REMARK 2.4.** *There exists a number of algorithms for the problem of network inference, and we will comment on two most relevant techniques. First, we note that the ARACNE algorithm [46] attempts to infer a (non-causal) interaction network based on mutual information. The ARACNE algorithm first computes the mutual information between all pairs of nodes/variables, filtering out the nonsignificant ones, and then enumerates through all triplets and removes links based on the data processing inequality. It was proven to correctly infer the undirected network under the assumptions that (i) mutual information are estimated without error, and (ii) the network is a tree [46]. Second, the PC algorithm developed by Spirtes, Glymour, and Scheines removes non-causal links by potentially testing all combinations of conditioning subsets, and was proven to correctly infer general causal networks if the conditional independence between the variables can be perfectly examined [66]. Runge et. al. [60, 61] recently utilized the PC algorithm to infer causal networks by establishing the conditional dependence/independence via estimation of appropriately defined conditional mutual information between time-lagged variables. We note that whereas we utilize Algorithm 2.2 for the divisive step in network inference, an alternative would be to utilize the PC algorithm for the divisive step. Although the accuracy versus efficiency tradeoff for such a modification has yet to be tested, we expect that it may be helpful specifically for inferring the causal parents for nodes with large degree, suggesting that in practical applications one may wish to switch back and forth between Algorithm 2.2 and the PC Algorithm for the divisive step, depending on a node's degree.*

**3. Application to Gaussian Process: Analytical Results.** In this section we make analytical comparison among three approaches to causal network inference:

---

<sup>6</sup>Numerically estimated causation entropy is always positive due to finite sample size and numerical precision. In practice, one needs to use a statistical test (e.g., permutation test as described in Section 4) to examine the conditions  $x > 0$  in Algorithm 2.1 and  $C_{j \rightarrow I|(K-\{j\})} = 0$  in Algorithm 2.2.

causation entropy, transfer entropy [63], and conditional Granger causality [29, 30]. The next section will be devoted to the exploration of the numerical properties of these approaches for general random networks.

While information-theoretic approaches including causation entropy do not require stringent model assumptions, a linear model must be assumed to offer a fair comparison with the conditional Granger causality. As a benchmark example, we focus on the following linear discrete stochastic network dynamics

$$(3.1) \quad X_t^{(i)} = \sum_{j \in N_i} A_{ij} X_{t-1}^{(j)} + \xi_t^{(i)} \quad (\text{or in matrix form: } X_t = AX_{t-1} + \xi_t).$$

Here  $X_t^{(i)} \in \mathbb{R}$  represents the state of node  $i$  at time  $t$  ( $i \in \{1, 2, \dots, n\}, t \in \mathbb{N}$ ),  $\xi_t^{(i)} \in \mathbb{R}$  represents noise, and  $A_{ij} X_{t-1}^{(j)}$  models the influence of node  $j$  on node  $i$ . Equation (3.1) finds application in a broad range of areas, including time series analysis (as a multivariate linear autoregressive process [10]), information theory (as a network communication channel [16]), and nonlinear dynamical systems (as a linearized stochastic perturbation around equilibrium states [43]). It is straightforward to check that Eq. (3.1) is a special case of the general network stochastic process, Eq. (2.4), and asymptotically (as  $t \rightarrow \infty$ ) satisfied the Markov assumptions in Eq. (2.8).

### 3.1. Analytical Properties of the Solution.

**3.1.1. Solution Formula.** Defining  $X_0 = \xi_0$  for convenience, the solution to Eq. (3.1) can be expressed as

$$(3.2) \quad X_t = \sum_{k=0}^t A^k \xi_{t-k}.$$

We assume that  $\xi_t^{(i)}$  are i.i.d Gaussian random variables with zero mean and finite nonzero variance, denoted as  $\xi_t^{(i)} \sim N(0, \sigma_i^2)$  with  $\sigma_i > 0$ . Therefore,

$$(3.3) \quad \xi_t \sim N(0, S),$$

where the covariance matrix  $S$  is defined by  $S_{ij} = \delta_{ij} \sigma_i^2$  with  $\delta$  denoting the Kronecker delta. It follows that

$$(3.4) \quad \begin{cases} \mathbb{E}[\xi_t^{(i)}] = 0, \\ \text{Cov}(\xi_t^{(i)}, \xi_\tau^{(j)}) = \delta_{ij} \delta_{t\tau}. \end{cases}$$

Note that a random variable obtained by an affine transformation of a Gaussian variable is also Gaussian. For example, if  $Y = [Y_1; Y_2]$  is Gaussian, the distribution of  $Y_1$  conditioned on  $Y_2$  is also Gaussian [20]. The proposition below follows by expressing random variables via appropriate affine transformations of  $\xi_t$ 's.

**PROPOSITION 3.1.** *Let  $I$  and  $K$  be any subsets of  $\mathcal{V}$ . Let  $t \in \mathbb{N}$  and  $\tau \in \{0\} \cup \mathbb{N}$ . The conditional distribution of  $X_{t+\tau}^{(I)}$  given  $X_t^{(K)}$  is Gaussian.*

**3.1.2. Covariance Matrix.** Under an affine transformation from Gaussian variable  $Y$  to  $Z$  as  $Z = CY + d$ , the mean and covariance of  $Y$  and  $Z$  are related by:  $\mu_Z = C\mu_Y + d$  and  $\Sigma_Z = C\Sigma_Y C^\top$  [20]. We consider covariance matrices  $\Phi(\tau, t)$ , where the  $(i, j)$ -th entry of  $\Phi(\tau, t)$  is defined as

$$(3.5) \quad \Phi(\tau, t)_{ij} \equiv \text{Cov}[x_{t+\tau}^{(i)}, x_t^{(j)}].$$

It follows from Eqs. (3.2) and (3.3) that

$$(3.6) \quad X_t \sim N(0, \Phi(0, t)), \text{ where } \Phi(0, t) = \sum_{k=0}^t A^k S (A^k)^\top.$$

In the following we prove a sufficient condition for the converge of the covariance matrix  $\Phi(0, t)$  as time  $t \rightarrow \infty$ . Denote the *spectral radius* of a square matrix  $M$  by

$$(3.7) \quad \rho_M \equiv \max\{|\lambda| : \lambda \text{ is an eigenvalue of } M\}.$$

Note that  $\rho_M = \rho_{M^\top}$  since a square matrix and its transpose have the same set of eigenvalues. For the dynamical system defined by Eq. (3.1), matrices  $A$  with  $|\rho_A| < 1$  are the only matrices for which the underlying system poses a stable equilibrium in the absence of noise. We refer to these matrices as stable.

DEFINITION 3.2 (Stable Matrix). *Matrix  $M$  is stable if  $\rho_M < 1$ .*

The following is a known result from classical matrix theory [36].

THEOREM 3.3 (Convergence of Matrix Series [36]). *The matrix series  $\sum_{k=0}^{\infty} M_k$  converges if the scalar series  $\sum_{k=0}^{\infty} \|M_k\|$  under any induced norm  $\|\cdot\|$  converges.*

Note that it is possible for the matrix series  $\sum_{k=0}^{\infty} M_k$  to be convergent while the corresponding scalar series  $\sum_{k=0}^{\infty} \|M_k\|$  diverges, analogous to the possibility of a scalar series that is convergent but not absolutely convergent. Next we state and prove a sufficient condition under which the matrix series in Eq. (3.6) converges.

PROPOSITION 3.4 (Convergence of the Covariance). *The series  $\sum_{k=0}^{\infty} A^k S (A^k)^\top$  converges if  $A$  is stable.*

*Proof.* Let  $\|\cdot\|$  be any induced norm. Then  $\|A^k S (A^k)^\top\| \leq \|A^k\| \cdot \|S\| \cdot \|(A^\top)^k\|$  for any  $k \in \mathbb{N}$ . Gelfand's formula (see Ref. [25]) implies that

$$(3.8) \quad \lim_{k \rightarrow \infty} \|A^k\|^{1/k} = \lim_{k \rightarrow \infty} \|(A^\top)^k\|^{1/k} = \rho_A.$$

On the other hand,  $\lim_{k \rightarrow \infty} \|S\|^{1/k} = 1$ . Therefore,

$$\lim_{k \rightarrow \infty} \|A^k S (A^k)^\top\|^{1/k} \leq \lim_{k \rightarrow \infty} (\|A^k\| \cdot \|S\| \cdot \|(A^\top)^k\|)^{1/k} = \rho_A^2 < 1,$$

where the last inequality follows from the fact that  $A$  is stable. Hence the scalar series  $\sum_{k=0}^{\infty} \|A^k S (A^k)^\top\|_2$  is convergent. The proposition follows by Theorem 3.3.  $\square$

For the remainder of this section, it will be assumed that  $A$  is stable in Eq. (3.1). As  $t \rightarrow \infty$ , we drop the second argument in  $\Phi(0, t)$  and define the *asymptotic covariance matrix*

$$(3.9) \quad \Phi(0) \equiv \lim_{t \rightarrow \infty} \Phi(0, t) = \sum_{k=0}^{\infty} A^k S (A^k)^\top.$$

It follows that  $\Phi(0)$  satisfies an algebraic equation given by the proposition below.

PROPOSITION 3.5 (Asymptotic Covariance Matrix). *Assume that  $A$  is stable. The asymptotic covariance matrix  $\Phi(0) = \sum_{k=0}^{\infty} A^k S (A^k)^\top$  satisfies the equation*

$$(3.10) \quad A\Phi(0)A^\top - \Phi(0) + S = 0.$$

*Proof.* Since  $A$  is stable, both of the two matrix series below converge:

$$\begin{cases} \Phi(0) = S + ASA^\top + A^2S(A^2)^\top + A^3S(A^3)^\top + \dots \\ A\Phi(0)A^\top = ASA^\top + A^2S(A^2)^\top + A^3S(A^3)^\top + \dots \end{cases}$$

Subtracting the two equations gives the result of the proposition.  $\square$

Equation (3.10) is a (discrete) Lyapunov equation which often appears in stability analysis and optimal control problems [59]. Using “ $\otimes$ ” as the Kronecker product and “vec” for the operation of transforming a square matrix to a column vector by stacking the columns of the underlying matrix in order, Eq. (3.10) can be converted into:

$$(3.11) \quad (I_{n^2} - A \otimes A) \text{vec}(\Phi(0)) = \text{vec}(S),$$

where  $I_{n^2}$  denotes the identity matrix of size  $n^2$ -by- $n^2$ . Matrix  $\Phi(0)$  can be computed by either solving Eq. (3.10) through iterative methods (see Ref. [5]) or by directly solving Eq. (3.11) as a linear system. In practice, we found the iterative approach to be numerically more efficient and stable compared to direct inversion.

Covariance matrices are in general positive semidefinite [20]. For for the network dynamics defined in Eq. (3.1), we show that they are indeed positive definite.

**PROPOSITION 3.6 (Positive Definiteness of the Covariance Matrix).** *The covariance matrix  $\Phi(0, t)$  is positive definite for any  $t \in \mathbb{N}$ . The asymptotic covariance matrix  $\Phi(0)$  is also positive definite.*

*Proof.* For any unit vector  $v \in \mathbb{R}^n$ ,  $v^\top A\Phi(0, 0)A^\top v = (A^\top v)^\top A^\top v \geq 0$ . From Eqs. (3.2) and (3.3), for any  $t \in \mathbb{N}$ ,  $\Phi(0, t) = A\Phi(0, t-1)A^\top + S$ . By induction,

$$(3.12) \quad \begin{aligned} v^\top \Phi(0, t)v &= v^\top A\Phi(0, t-1)A^\top v + v^\top Sv \\ &\geq (A^\top v)^\top \Phi(0, t-1) (A^\top v) + \min_i \sigma_i^2 \geq \min_i \sigma_i^2 > 0. \end{aligned}$$

This shows that  $\Phi(0, t)$  is positive definite (indeed we have:  $\rho_{\Phi(0, t)} \geq \min_i \sigma_i^2 > 0$ ). Taking  $t \rightarrow \infty$  in the above estimate also shows that  $\Phi(0)$  is positive definite.  $\square$

**3.1.3. Time-Shifted Covariance Matrices.** We define the time-shifted covariance matrix  $\Phi(\tau, t)$  for each  $t \in \mathbb{N}$  (time) and  $\tau \in \mathbb{N}$  (positive time shift between states). If  $A$  is stable, then the covariance matrix  $\Phi(t, \tau)$  converges for each time shift  $\tau$  as  $t \rightarrow \infty$ . The (asymptotic) covariance matrices with different time shifts are related by a simple algebraic equation given in the following proposition.

**PROPOSITION 3.7 (Relationship Between Time-Shifted Covariance Matrices).** *Assume that  $A$  is stable. For each  $\tau \in \mathbb{N}$ , the following limit exists*

$$\lim_{t \rightarrow \infty} \Phi(\tau, t) = \Phi(\tau),$$

where matrix  $\Phi(\tau)$  satisfies

$$(3.13) \quad \Phi(\tau) = A\Phi(\tau-1) = A^2\Phi(\tau-2) = \dots = A^\tau\Phi(0).$$

*Proof.* For every  $\tau \in \mathbb{N}$  and  $t \in \mathbb{N}$ , it follows that

$$(3.14) \quad \Phi(\tau, t)_{ij} = \mathbb{E} \left[ \sum_{k=1}^n a_{ik} x_{t+\tau-1}^{(k)} + \xi_{t+\tau}^{(i)}, x_t^{(j)} \right] = \sum_{k=1}^n a_{ik} \Phi(\tau-1, t)_{kj}.$$

Therefore, the matrix  $\Phi(\tau, t)$  satisfies

$$(3.15) \quad \Phi(\tau, t) = A\Phi(\tau-1, t) = A^2\Phi(\tau-2, t) = \dots = A^\tau\Phi(0, t).$$

Taking the limit as  $t \rightarrow \infty$  in and making use of the fact that  $A$  is stable, we reach the conclusion of the proposition.  $\square$



**3.2. Analytical Expressions of Causation Entropy.** Here we provide analytical expressions for causation entropy of the Gaussian process described in Eq. (3.1). Because causation entropy can be interpreted as a generalization of both transfer entropy and conditional Granger causality under the appropriate selection of nodes  $i$  and  $j$  and the conditioning set  $K$ , these results also provide analytical expressions for transfer entropy and conditional Granger causality.

**3.2.1. Joint entropy expressions.** Let  $\Sigma$  be the covariance matrix of a multivariate Gaussian variable  $X \in \mathbb{R}^n$  (i.e.,  $X \sim N(\boldsymbol{\mu}, \Sigma)$ ), it follows that [1]

$$(3.16) \quad h(X) = \frac{1}{2} \log[\det(\Sigma)] + \frac{1}{2} n \log(2\pi e).$$

Note that the right hand side of the above is actually an upper bound for a general random variable (i.e., the equality “=” becomes inequality “ $\leq$ ” [16]). Therefore, a Gaussian variable maximizes entropy among all variables of equal covariance.

The random variable  $X_t$  is Gaussian and converges to  $N(0, \Phi(0))$  as  $t \rightarrow \infty$ . For an arbitrary subset of the nodes  $K = \{k_1, k_2, \dots, k_\ell\}$ . The joint entropy is

$$(3.17) \quad h(X^{(K)}) = \lim_{t \rightarrow \infty} h(X_t^{(K)}) = \frac{1}{2} \log(|\Phi_{KK}(0)|) + \log(2\pi e).$$

Here we have introduced the notation

$$(3.18) \quad \Phi_{IJ}(0) \equiv P(I)\Phi(0)P(J)^\top,$$

where for a set  $K = \{k_1, k_2, \dots, k_\ell\}$ ,  $P(K)$  is the  $\ell$ -by- $n$  projection matrix defined as

$$(3.19) \quad P(K)_{ij} = \delta_{k_i, i}$$

**3.2.2. Causation Entropy.** For the Gaussian process given by Eq. (3.1), we obtain the analytical expression of causation entropy as

$$(3.20) \quad C_{J \rightarrow I|K} = \frac{1}{2} \log \left( \frac{\det [\Phi(0)_{II} - \Phi(1)_{IK} \Phi(0)_{KK}^{-1} \Phi(1)_{IK}^\top]}{\det [\Phi(0)_{II} - \Phi(1)_{I, K \cup J} \Phi(0)_{K \cup J, K \cup J}^{-1} \Phi(1)_{I, K \cup J}^\top]} \right)$$

If  $J = \{j\}$  and  $I = \{i\}$ , this equation simplifies to

$$(3.21) \quad C_{j \rightarrow i|K} = \frac{1}{2} \log \left( \frac{\Phi(0)_{ii} - \Phi(1)_{iK} \Phi(0)_{KK}^{-1} \Phi(1)_{iK}^\top}{\Phi(0)_{ii} - \Phi(1)_{i, K \cup \{j\}} \Phi(0)_{K \cup \{j\}, K \cup \{j\}}^{-1} \Phi(1)_{i, K \cup \{j\}}^\top} \right).$$

**3.2.3. Transfer Entropy.** Recall that causation entropy recovers transfer entropy when  $K = \{i\}$ . Letting  $K = \{i\}$  in the formula above gives the transfer entropy (with single time lag) for multivariate Gaussian variables:

$$(3.22) \quad T_{j \rightarrow i} = C_{j \rightarrow i|i} = \frac{1}{2} \log \left( 1 + \frac{\alpha_{ij}}{\beta_{ij} - \alpha_{ij}} \right),$$

where  $\begin{cases} \alpha_{ij} \equiv (\Phi(0)_{ii} \Phi(1)_{ij} - \Phi(0)_{ij} \Phi(1)_{ii})^2, \\ \beta_{ij} \equiv (\Phi(0)_{ii}^2 - \Phi(1)_{ii}^2) (\Phi(0)_{ii} \Phi(0)_{jj} - \Phi(0)_{ij}^2). \end{cases}$

It follows that  $\beta_{ij} \geq \alpha_{ij} \geq 0$ , and therefore  $T_{j \rightarrow i} \geq 0$  ( $T_{i \rightarrow i} = 0$ ). Furthermore,

$$(3.23) \quad T_{j \rightarrow i} = 0 \iff \alpha_{ij} = 0 \iff \sum_{k=1}^n A_{ik} (\Phi(0)_{ii} \Phi(0)_{kj} - \Phi(0)_{ij} \Phi(0)_{ki}) = 0.$$

**3.2.4. Conditional Granger Causality.** As shown in Ref. [3], when the random variables are Gaussian, expression of Granger Causality is equivalent as that of transfer entropy (and also causation entropy introduced here). In fact, for Gaussian variables, the Granger Causality from  $j$  to  $i$  without conditioning equals  $2C_{j \rightarrow i}$ , while the conditional Granger causality (with full conditioning) equals  $2C_{j \rightarrow i | (\mathcal{V} - \{j\})}$ .

**3.3. Analytical Results for Directed Linear Chain, Directed Loop, and Directed Trees.** We derive expressions of transfer entropy and causation entropy for several classes of networks including directed linear chains, directed loops, and directed trees. These results highlight that although transfer entropy may indicate the direction of information flow between two nodes, its application to causal network inference is often unjustified as it cannot distinguish between direct and indirect causal relationships (unless appropriate conditioning is adopted as in causation entropy).

**3.3.1. Directed Linear Chain.** Denote a directed linear chain of  $n$  nodes as

$$(3.24) \quad 1 \rightarrow 2 \rightarrow 3 \cdots \rightarrow n.$$

For simplicity we assume that all links have the same weight  $w = 1$ . Consequently, the corresponding adjacency matrix  $A = [A_{ij}]_{n \times n}$  is given by

$$(3.25) \quad A_{ij} = \delta_{i,j+1}.$$

It follows that  $\rho_A = 0$  and therefore  $A$  is stable. By inverting the lower-triangular matrix  $(I_{n^2} - A \otimes A)$  in Eq. (3.11) and applying Eq. (3.13), we obtain that

$$(3.26) \quad \begin{cases} \Phi(0)_{ij} = \delta_{ij} \sum_{k=1}^j \sigma_k^2, \\ \Phi(1)_{ij} = \delta_{i,j+1} \sum_{k=1}^j \sigma_k^2. \end{cases}$$

Letting  $K = \emptyset$  and  $K = \{i\}$  respectively in Eqs. (3.21) and (3.22), it follows that

$$(3.27) \quad C_{j \rightarrow i} = T_{j \rightarrow i} = \frac{1}{2} \delta_{i,j+1} \log \left( 1 + \frac{\sum_{k=1}^j \sigma_k^2}{\sigma_i^2} \right).$$

Therefore, for the directed linear chain defined in Eq. (3.25), transfer entropy  $T_{j \rightarrow i} = C_{j \rightarrow i}$ , and it is positive if and only if there is a direct link  $j \rightarrow i$ , i.e.,

$$(3.28) \quad C_{j \rightarrow i} = T_{j \rightarrow i} > 0 \Leftrightarrow A_{ij} = 1, \text{ and } C_{j \rightarrow i} = T_{j \rightarrow i} = 0 \Leftrightarrow A_{ij} = 0.$$

Interestingly, both causation entropy  $C_{j \rightarrow j+1}$  and transfer entropy  $T_{j \rightarrow j+1}$  increase monotonically as a function of  $j$ , and the values only depend on part of the chain from the top node (node 1) to node  $j+1$  and not on the rest of the network. Interpreting the monotonicity in term of the network structure, the closer node  $j$  is to the end of the chain, effectively the more information is transferred through the directed link  $j \rightarrow j+1$ . Figure. 3.1(a) illustrates this via a network of  $n = 1000$  nodes.

**3.3.2. Directed Loop.** Consider now a directed loop with  $n$  nodes, denoted as

$$(3.29) \quad 1 \rightarrow 2 \rightarrow 3 \cdots \rightarrow n \rightarrow 1.$$

Let  $w > 0$  be the uniform link weight. It follows that  $\rho_A = w$ . Thus, for the adjacency matrix  $A$  to be stable, we must have  $w < 1$ . To keep the symmetry of the problem, we further assume that the variance of noise is the same at each node, therefore

$$(3.30) \quad \sigma^2 \equiv \sigma_1^2 = \sigma_2^2 = \dots \sigma_n^2.$$

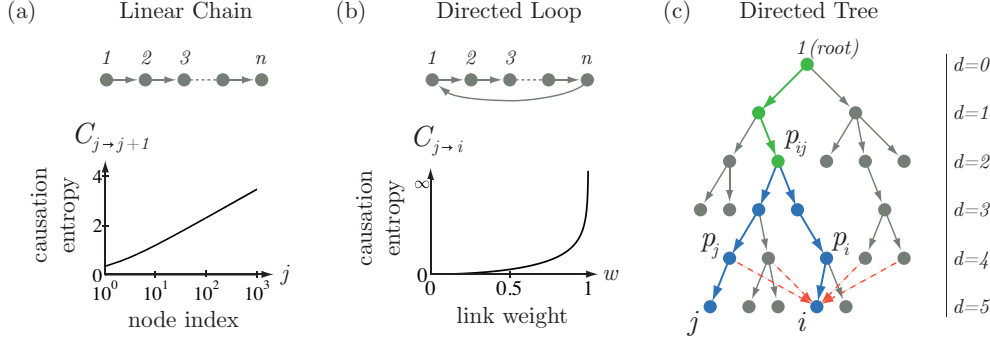


FIG. 3.1. *Causation Entropy and transfer entropy for a Gaussian process on three classes of networks.* (a) For directed linear chains, both causation entropy and transfer entropy correctly identify the network as  $C_{j \rightarrow i} = T_{j \rightarrow i} > 0$  iff  $i = j+1$  (otherwise  $C_{j \rightarrow i} = T_{j \rightarrow i} = 0$ ). The dependence of  $C_{j \rightarrow j+1}$  on node index  $j$  is given by Eq. (3.27) and plotted. (b) For directed loops, causation entropy and transfer entropy again correctly identify the network topology with  $C_{j \rightarrow i} = T_{j \rightarrow i} > 0$  iff  $j \rightarrow i$ . The dependence of  $C_{j \rightarrow i}$  on link weight  $w$  is given by Eq. (3.33) as shown. (c) For directed trees, causation entropy given by Eq. (3.41) correctly identifies the network topology based on Eq. (3.43). In contrast, transfer entropy without appropriate conditioning infers many links that do not exist in the actual network (red dashed lines), as described by Eq. (3.42).

The entries in  $\Phi(0, t)$  satisfy

$$(3.31) \quad \Phi(0, t)_{ij} = w^2 \Phi(0, t-1)_{p_i, p_j} + \delta_{ij} \sigma^2,$$

where  $p_i$  denotes the unique node that directly links to node  $i$ . Taking the limit as  $t \rightarrow \infty$  and solve the resulting recursive equations, we obtain that for

$$(3.32) \quad \begin{cases} \Phi(0)_{ij} = \delta_{ij} \sigma^2 / (1 - w^2), \\ \Phi(1)_{ij} = \delta_{p_i, j} \sigma^2 w / (1 - w^2). \end{cases}$$

where the second equation is obtained through  $\Phi(0)_{ij}$  and Eq. (3.13). Letting  $K = \emptyset$  and  $K = \{i\}$  respectively in Eqs. (3.21) and (3.22), we conclude that

$$(3.33) \quad C_{j \rightarrow i} = T_{j \rightarrow i} = \frac{1}{2} \delta_{p_i, j} \log \left( \frac{1}{1 - w^2} \right).$$

Note that causation entropy and transfer entropy equal and do not depend on the noise variation  $\sigma^2$ , and they are positive if and only if there is a direct link  $j \rightarrow i$ , i.e.,

$$(3.34) \quad C_{j \rightarrow i} = T_{j \rightarrow i} > 0 \Leftrightarrow A_{ij} = 1, \quad \text{and} \quad C_{j \rightarrow i} = T_{j \rightarrow i} = 0 \Leftrightarrow A_{ij} = 0.$$

By symmetry, causation entropy and transfer entropy through each directed link is the same. As the link weight  $w$  increases in  $(0, 1)$ , both increase monotonically in  $(0, \infty)$ . The larger the link weight  $w$  is, the larger amount of information is transferred via each directed link, as intuitively expected. Also see Fig. 3.1(b) as an illustration.

**3.3.3. Directed Trees.** We now consider directed tree networks with uniform link weight  $w = 1$  and unit node variance<sup>7</sup>

$$(3.35) \quad \sigma_1^2 = \sigma_2^2 = \dots \sigma_n^2 = 1.$$

<sup>7</sup>Similar results hold for trees with general link weights and node variances but the corresponding equations are too cumbersome to list.

A directed tree has one root (indexed as node 1 without loss of generality) and each non-root node  $i$  ( $i \neq 1$ ) has exactly one *ancestor*, denoted by  $p_i$ . The corresponding adjacency matrix  $A = [A_{ij}]_{n \times n}$  thus satisfies

$$(3.36) \quad A_{ij} = (1 - \delta_{i1})\delta_{i,p_i}.$$

It can be shown that  $\rho_A = 0$ . For  $i \neq 1$ , we denote the directed path from 1 to  $i$  by

$$(3.37) \quad 1 = p_i^{(d_i)} \rightarrow p_i^{(d_i-1)} \rightarrow \dots \rightarrow p_i^{(1)} \equiv p_i \rightarrow p_i^{(0)} \equiv i,$$

where  $d_i$  is the *depth* of node  $i$  in the tree (for node 1, we define its depth  $d_1 = 0$ ). Thus, the highest node in the tree is the root, and the lowest nodes have the greatest depth. For any two nodes  $(i, j)$ , we denote their *lowest common ancestor* by  $p_{ij}$ , i.e.,

$$(3.38) \quad p_{ij} = \arg \max_{\{k | \exists \ell, m \geq 0, s.t., p_i^{(\ell)} = p_j^{(m)}\}} d_k.$$

The covariance matrix  $\Phi(0, t)$  satisfies

$$(3.39) \quad \Phi(0)_{ij} = \delta_{1i}\delta_{1j}\sigma_1^2 + (1 - \delta_{1i})(1 - \delta_{1j})[\Phi(0)_{p_i,p_j} + \delta_{ij}].$$

We solve these recursive equations to obtain

$$(3.40) \quad \begin{cases} \Phi(0)_{ij} = \delta_{d_i,d_j}(d_{p_{ij}} + 1) \\ \Phi(1)_{ij} = (1 - \delta_{i1})\delta_{d_i,d_j+1}(d_{p_{ij}} + 1), \end{cases}$$

where  $p_{ij}$  is defined in Eq. (3.38) and  $\Phi(1)_{ij}$  is obtained by  $\Phi(1) = A\Phi(0)$ .

We calculate causation entropy and transfer entropy through Eqs. (3.21) and (3.22):

$$(3.41) \quad C_{j \rightarrow i} = T_{j \rightarrow i} = \frac{1}{2}\delta_{d_i,d_j+1} \log \frac{(d_i + 1)(d_j + 1)}{(d_i + 1)(d_j + 1) - (d_{p_{ij}} + 1)^2}.$$

Note that in general  $0 \leq d_{p_{ij}} \leq \min\{d_i, d_j\}$ . Thus  $C_{j \rightarrow i} = T_{j \rightarrow i} \leq \frac{1}{2} \log(1 + d_i)$ , with equality if and only if  $j$  is the ancestor of  $i$  (i.e.,  $j = p_i = p_{ij}$ ). Therefore, we have

$$(3.42) \quad \begin{cases} T_{j \rightarrow i} > 0 \Leftrightarrow d_i = d_j + 1 \Leftarrow A_{ij} = 1 \text{ (but } T_{j \rightarrow i} > 0 \not\Leftarrow A_{ij} = 1); \\ T_{j \rightarrow i} = 0 \Leftrightarrow d_i \neq d_j + 1 \Rightarrow A_{ij} = 0 \text{ (but } A_{ij} = 0 \not\Rightarrow T_{j \rightarrow i} = 0). \end{cases}$$

In other words, transfer entropy being positive (without appropriate conditioning) corresponds to a superset of the links that actual exist in a directed tree, and the inferred network using this criterion will potentially contain many false positives. See Fig. 3.1(c) as an example. On the other hand, for a given node  $i \neq 1$ , we have

$$(3.43) \quad \begin{cases} p_i = \arg \max_j C_{j \rightarrow i}, \\ C_{j \rightarrow i | \{p_i\}} = 0. \end{cases}$$

Therefore, for each node  $i$ , the node  $j$  that maximizes causation entropy  $C_{j \rightarrow i}$  among all nodes is inferred as the causal parent of  $i$ . Conditioned on this node, the causation entropy from any other node to  $i$  will become zero, indicating no other directed links to node  $i$ . This causation entropy based procedure allows for exact and correct inference of the underlying causal network, a directed tree.

**4. Application to Gaussian Process: Numerical Results.** In this section, we illustrate that causal network inference by optimal causation entropy is reliable and efficient for the Gaussian process, Eq. (3.1), on large random networks.

**4.1. Random Network Model and Time Series Generation.** We consider signed Erdős-Rényi networks, which is a generation of its original model [8]. In particular, each network consists of  $n$  nodes ( $\mathcal{V} = \{1, 2, \dots, n\}$ ), such that each directed link  $j \rightarrow i$  is formed independently with equal probability  $p$ , giving rise to a directed network with approximately  $n^2p$  directed links. For generality, we allow the link weight of each link  $j \rightarrow i$  to be either positive ( $A_{ij} = w$ ) or negative ( $A_{ij} = -w$ ), with equal probability. Recalling that the network adjacency matrix  $A$  is defined entry-wise by  $A_{ij} \in \{w, -w\}$  iff there exists a directed link  $j \rightarrow i$  (otherwise  $A_{ij} = 0$ ), the link weight  $w$  may be selected to tune the spectral radius  $\rho(A)$  of matrix  $A$ .

We generate time series from the stochastic equation, Eq. (3.1), where matrix  $A$  is obtained from the network model and random variables  $\xi_t \sim \mathcal{N}(0, S)$ , where the covariance matrix  $S$  is taken to be the identity matrix of size  $n \times n$ . To reduce transient effects, for a given sample size  $T$  we solve Eq. (3.1) for  $10T$  time steps and only use the final 10% of the resulting time series.

To summarize, our numerical experiments contain parameters:  $n$  (network size),  $p$  (connection probability),  $\rho(A)$  (spectral radius of  $A$ ), and  $T$  (sample size).

**4.2. Practical Considerations for Network Inference.** We have established by Theorems 2.2 and 2.3 and Lemmas 2.4 and 2.5 that in theory, exact network inference can be achieved by optimal causation entropy, which involves implementing Algorithms 2.1 (Aggregative Discovery) and 2.2 (Progressive Removal) to correctly identify the set of causal parents  $N_i$  for each node  $i \in \mathcal{V}$ .

In practice, the success of our optimal causation entropy approach (and in fact, any entropy-based approaches) depends crucially on reliable estimation of the relevant entropies in question from data. This leads to two practical challenges.

(1) Entropies must be *estimated* from *finite* time series data. While there are several techniques for estimating entropies for general multivariate data, the accuracy of such estimations are increasingly inaccurate for small sample sizes and high-dimensional random variables [52]. In this research, we side-step this computational complexity by using knowledge of the asymptotic functional form for the entropy of the Gaussian Process, where the covariance matrices  $\Phi(0)$  and  $\Phi(1)$  in Eqs. (3.20) and (3.21) are estimated directly from the time series data.

(2) Application of the theoretical results rely on determining whether the causation entropy  $C_{j \rightarrow i|K} > 0$  or  $C_{j \rightarrow i|K} = 0$ . However, the estimated value of  $C_{j \rightarrow i|K}$  based on sample covariances is necessarily positive given finite sample size and finite numerical precision. Therefore, a statistical test must be used to assess the significance of the observed positive causation entropy. We here adopt a widely used approach in non-parametric statistics, called the *permutation test*<sup>8</sup>. Specifically, we propose the following permutation test based on the null hypothesis that causation entropy  $C_{j \rightarrow i|K} = 0$ : first perform  $r$  random (temporal) permutations of the time series  $\{X_t^{(j)}\}$ , leaving the rest of the data unchanged; we then construct an empirical cumulative distribution  $\hat{F}(x)$  of the estimated causation entropy from the permuted

<sup>8</sup>The idea of a permutation test is to perform (large number of) random permutations of a subset of the data leaving the rest unchanged, giving rise to an empirical distribution of the static of interest. The observed statistic from the original data is then located on this empirical distribution in order to associate its statistical significance [28].

time series<sup>9</sup>; finally, given a prescribed significance level  $\theta$ , the observed  $C_{j \rightarrow i|K} = c$  is declared *significant* (i.e., the null hypothesis is rejected at level  $\theta$ ) if  $\hat{F}(c) > \theta$ .

To summarize, the inference algorithms contain two parameters to be used in the permutation test:  $r$  (number of random permutations) and  $\theta$  (significance threshold).

**4.3. Comparing Optimal Causation Entropy, Conditional Granger, and transfer entropy.** Here we compare the performance of three approaches of causal network inference: conditional Granger (see for example Ref. [23, 32]), transfer entropy (see Ref. [76] and the references therein), and optimal causation entropy (oCSE). In particular, the conditional Granger and transfer entropy approaches under consideration both estimate the entropy  $C_{j \rightarrow i|K}$  for each pair of nodes  $(i, j)$  independently, with the choice of  $K = \mathcal{V} - \{j\}$  in the case of conditional Granger and  $K = \{i\}$  in the case of transfer entropy. In both approaches, a causal link  $j \rightarrow i$  is inferred if the observed  $C_{j \rightarrow i|K} > 0$  is assessed as significant under the permutation test. The oCSE approach combines Algorithms 2.1 and 2.2 and the permutation test is used once per each iteration (line 2 of both algorithms).

The performance of the three approaches are quantified by two types of inference error: false negative ratio, denoted as  $\varepsilon_-$  and defined as the fraction of links in the original network that are not inferred; and false positive ratio, denoted as  $\varepsilon_+$  and defined as the fraction of non-existing links in the original networks that are inferred. In terms of the adjacency matrix  $A$  of the original network and that of the inferred network  $\hat{A}$ , these ratios can be computed as

$$(4.1) \quad \begin{cases} \varepsilon_- \equiv \frac{\text{number of } (i, j) \text{ pairs with } \chi_0(A)_{ij} = 1 \text{ and } \chi_0(\hat{A})_{ij} = 0}{\text{number of } (i, j) \text{ pairs with } \chi_0(A)_{ij} = 1}, \\ \varepsilon_+ \equiv \frac{\text{number of } (i, j) \text{ pairs with } \chi_0(A)_{ij} = 0 \text{ and } \chi_0(\hat{A})_{ij} = 1}{\text{number of } (i, j) \text{ pairs with } \chi_0(A)_{ij} = 0}. \end{cases}$$

For the random networks considered here, we found that the Algorithm 2.1 achieves almost the same accuracy as the combination of Algorithms 2.1 and 2.2. We therefore present results which are based on the numerical application of Algorithm 2.1 alone, leaving detailed numerical study of Algorithm 2.2 to future work.

Figure 4.1(a-b) shows that although the conditional Granger approach is theoretically correct and works well for small network size with sufficient samples, it suffers from increasing inference error as the network size increases and become extremely inaccurate when the network size  $n$  starts to surpass the sample size  $T$ . Such limitation is overcome by the oCSE approach, where both the false positive and false negative ratios remain close to zero as the network size increases. The reason that oCSE is accurate even as  $n$  increases is that it builds the causal parent set in an *aggregative* manner, therefore relying only on estimating entropy in relatively low dimensions (roughly the same dimension as the number of causal parents per node). In sharp contrast, the conditional Granger approach requires the estimation of entropy in the full  $n$ -dimensional space and therefore requires many (potentially exponentially) more samples to achieve the same accuracy when  $n$  becomes large.

Figure 4.1(c-d) shows that even for a sufficient number of samples, the transfer entropy approach without appropriate conditioning can lead to considerable inference error, and is therefore inherently unsound for causal network inference. In particular,

<sup>9</sup>The accuracy of this empirical distribution and therefore the permutation test increases with increasing number of permutations  $r$ . However, as  $r$  increases, the computational complexity also increases, scaling roughly as a linear function of  $r$ .

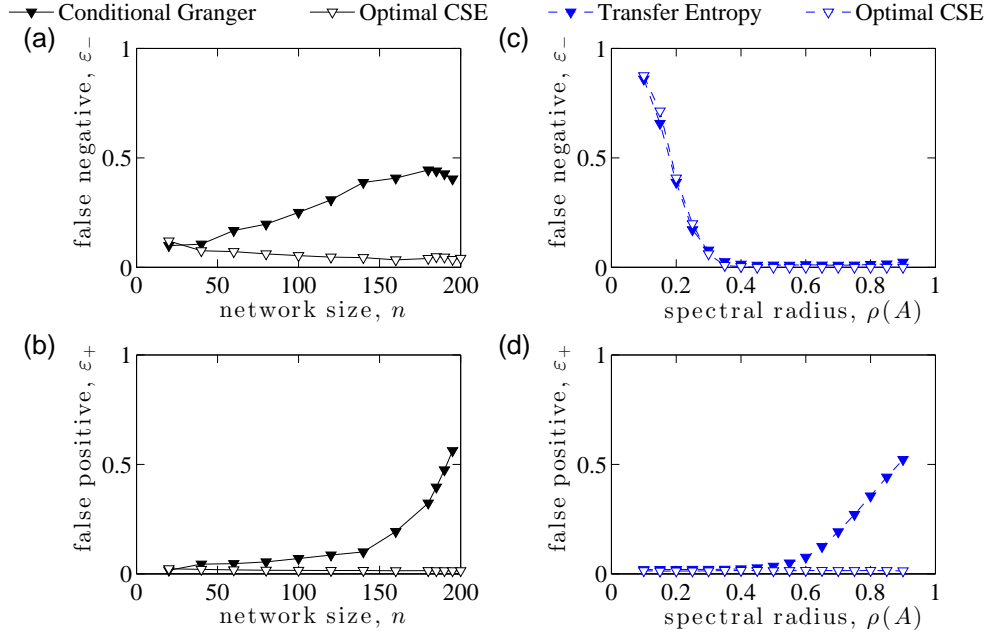


FIG. 4.1. Comparison of causal network inference approaches: conditional Granger, transfer entropy, and oCSE. The time series are generated from the Gaussian process defined in Eq. (3.1) using signed Erdős-Rényi networks (see Sec 4.2 for details). Two types of inference error are examined: false negative and false positive ratios, defined in Eq. (4.1). (a-b) Inference error as a function of network size  $n$  using conditional Granger versus oCSE approaches. Here the networks have fixed average degree  $np = 10$  and spectral radius  $\rho(A) = 0.8$ . Sample size is  $T = 200$ . (c-d) Inference error as a function of the spectral radius  $\rho(A)$  using transfer entropy versus oCSE approaches. Here the networks have fixed number of nodes  $n = 200$  and average degree  $np = 10$ . Sample size is  $T = 2000$ . For all three approaches we apply the permutation test using  $r = 100$  permutations and significance level  $\theta = 99\%$ . Each data point is obtained from averaging over 20 independent simulations of the network dynamics, Eq. (3.1).

although inference by both transfer entropy and oCSE give similar false negatives in the regime of  $\rho(A) \approx 0$  where the dynamics is dominated by noise and not the causal dependences, transfer entropy yields increasing false positives when the causal links dominate,  $\rho(A) \rightarrow 1$ . This is mainly due to the fact that as  $\rho(A) \rightarrow 1$ , indirect causal nodes become increasingly difficult to distinguish from direct ones without appropriate conditioning [71]. oCSE, on the other hand, consistently yields nearly zero false positive ratios in the entire range of  $\rho(A)$ . Interestingly, the spectral radius  $\rho(A)$  can be interpreted as the *information diffusion rate* on networks and found to be very close to criticality (i.e.,  $\rho(A) \approx 1$ ) in neuronal networks [39, 42].

These numerical experiments highlight that whereas the conditional Granger approach is inaccurate for  $T \lesssim n$  and the transfer entropy approach is inaccurate when  $\rho(A) \lesssim 1$ , the proposed oCSE approach overcomes both limitations and yields almost exact network inference even for limited sample size.

**4.4. Performance of Optimal Causation Entropy Approach for Causal Network Inference.** Having established the advantages of the oCSE approach, we now examine its performance under various parameter settings.

First, we examine the effect of the significance level  $\theta$  on the inference error. As

shown in Fig. 4.2(a-b), the false negative ratio  $\epsilon_-$  does not seem to depend on  $\theta$  and converges to zero as sample size  $T$  increases. On the other hand, as  $T \rightarrow \infty$ , the false positive ratio saturates at the level  $\epsilon_+ \sim (1 - \theta)$ , which is consistent with the implementation of the permutation test which rejects the null hypothesis at  $\theta$ . This observation suggests that in order to achieve higher accuracy given sufficient sample size, one should choose  $\theta$  as close to one as possible. The tradeoff in practice is that reliable implementation using larger  $\theta$  requires an increasing number of permutations and therefore increases the computational complexity of the inference algorithms.

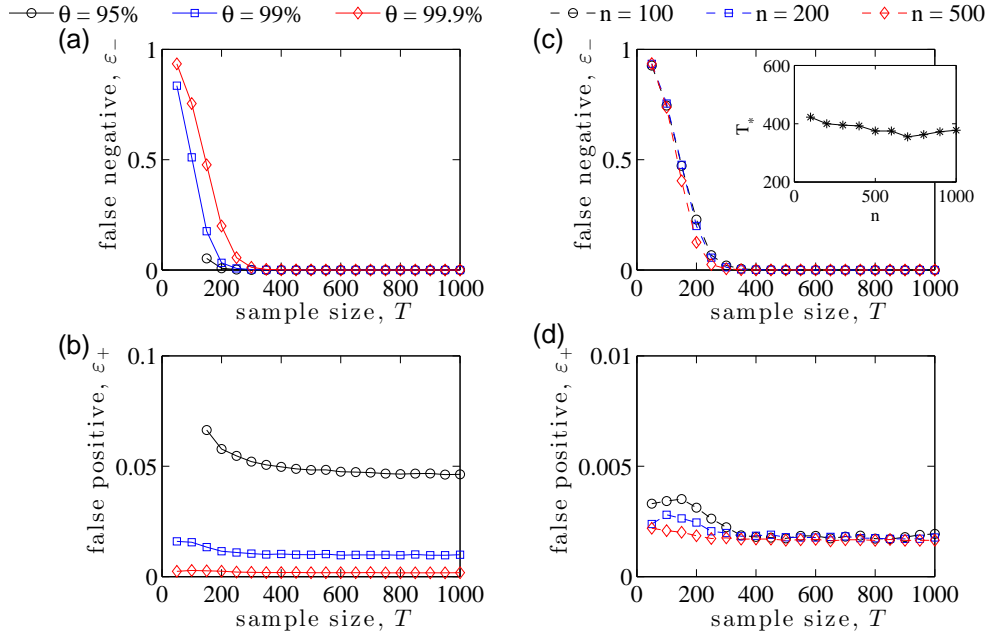


FIG. 4.2. Performance of the oCSE approach for causal network inference with different significance threshold for networks of various sizes. The time series are generated from the Gaussian process defined in Eq. (3.1) using signed Erdős-Rényi networks (see Sec 4.2 for details). False negative ratio (upper row) and false positive ratio (lower row) are defined in Eq. (4.1). (a-b) Inference error as a function of sample size  $T$  for various significance levels  $\theta$  used in the permutation test. Here networks have  $n = 200$  nodes with expected average degree  $np = 10$  and information diffusion rate  $\rho(A) = 0.8$ . (c-d) Inference error as a function of sample size  $T$  for various network sizes. Here networks have the same expected average degree  $np = 10$  and information diffusion rate  $\rho(A) = 0.8$ , and we use  $r = 1000$  permutations in the permutation test with  $\theta = 0.999$ . Note that all three false negative curves in (c) appear to converge for  $T \approx 300$ . The critical sample size  $T_*$  (defined as the minimum  $T$  for which  $\epsilon_- < 1 - \theta$ ) as a function of the network size  $n$  is shown in the inset of (c), suggesting the absence of scaling of  $T_*$  in terms of  $n$ . Each data point is obtained from averaging over 20 independent simulations of the network dynamics, Eq. (3.1).

Next, we investigate the effect of sample size  $T$  on the inference error for networks of different sizes. The results are shown in Fig. 4.2(c-d). As expected, when  $T$  increases, the false negative ratio decreases towards zero. Somewhat unexpectedly, the false positive ratio stays close to zero (in fact, close to the significance level  $\theta$ ) even for relatively small sample size ( $T$  as small as 50 for networks of up to 500 nodes). Furthermore, it appears that for networks of different sizes but the same average degree and information diffusion rate, the false negative ratios drop close to zero almost at the same sample size. To better quantify these effects, we define



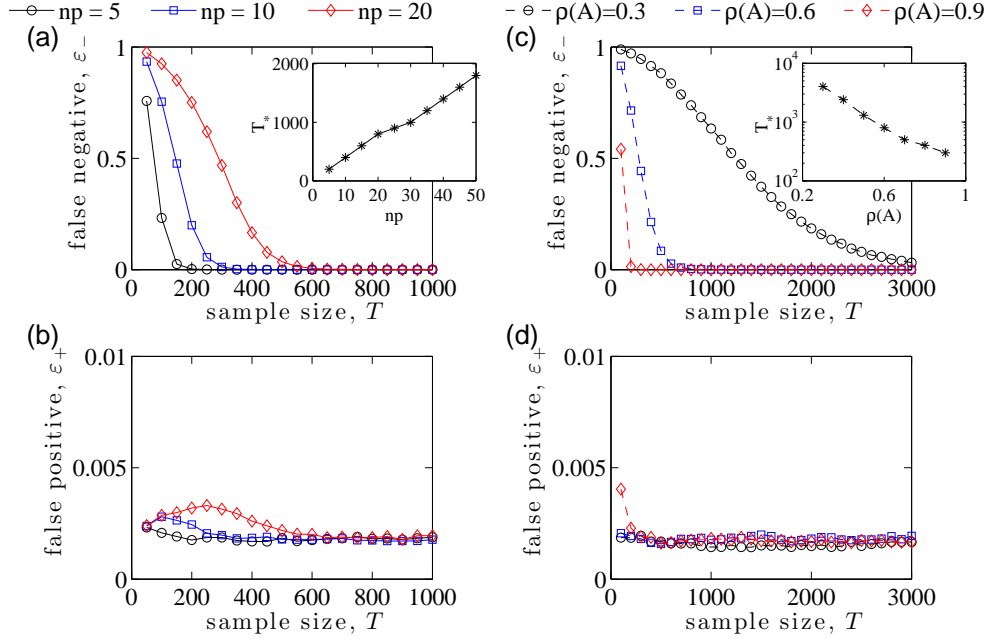


FIG. 4.3. Performance of the oCSE approach for causal network inference for networks with different average degree and spectral radius. The time series are generated from the Gaussian process defined in Eq. (3.1) using signed Erdős-Rényi networks (see Sec 4.2 for details). False negative ratio (upper row) and false positive ratio (lower row) are defined in Eq. (4.1). (a-b) Inference error as a function of sample size for networks with various average degree  $np$ . Here the networks have the same size  $n = 200$  and spectral radius  $\rho(A) = 0.8$ . The inset shows the critical sample size  $T_*$  (see text) as a function of  $np$ . (c-d) Inference error as a function of sample size for networks with various spectral radii  $\rho(A)$ . Here the networks have the same size  $n = 200$  and average degree  $n = 10$ . The permutation test used for the data in all panels involve  $r = 1000$  permutations with the significance threshold  $\theta = 0.999$ . Each data point is obtained from averaging over 20 independent simulations of the network dynamics, Eq. (3.1).

the critical sample size  $T_*$  as the smallest number of samples for which the false negative ratio falls below  $1 - \theta$ . As shown in the inset of Fig. 4.2(c), for networks with the same average degree and information diffusion rate, the critical sample size  $T_*$  remains mostly constant despite the increase of the network size. This result is unexpected. Traditionally, the network size  $n$  represents a lower bound on sample size  $T$  as any covariance matrix (e.g., application of the conditional Granger requires that  $T > n$  for the invertibility of the covariance matrices). Our result surprisingly indicates that sample size  $T$  does not need to scale with network size  $n$  for accurate network inference, and highlights the fact that the oCSE approach is *scalable* and *data efficient*, with accuracy depending *not* on the size of the network, but rather on other network characteristics such as the density of links and spectral radius.

To strengthen our claim that for Erdős-Rényi networks, performance of the causal inference by the oCSE approach depends on the density of links as measured by average degree and information diffusion rate as measured by the spectral radius rather than network size, we further investigate the dependence of inference error on these two additional parameters,  $np$  and  $\rho(A)$ . As shown in Fig. 4.3(a), for networks of the same size  $n = 200$  with fixed  $\rho(A) = 0.8$ , the larger the average degree  $np$ , the larger the number of samples required to reduce the false negative ratio to zero.

In fact, as shown in the inset of Fig. 4.3(a), the critical sample  $T_*$  to reach  $\varepsilon_- < 1 - \theta$  appears to scale *linearly* as a function of the average degree  $np$ , but not the network size (see the inset of Fig. 4.2(c)). On the other hand, Fig. 4.3(c-d) shows that the information diffusion rate,  $\rho(A)$ , seems to pose a harder constraint on accurate network inference: the smaller it is, the more samples that are needed for accuracy. In particular, as shown in the inset of Fig. 4.3(c), the critical sample size appears to increase *exponentially* as  $\rho(A)$  decreases towards zero. Interestingly, as shown in Fig. 4.3(b,d), the false positive ratios in both cases remain close to its saturation level around  $1 - \theta = 10^{-3}$  even for very small sample size ( $T \sim 50$ ), and this holds across networks with different average degree and different size (also see Fig. 4.2(d)).

To briefly summarize these numerical experiments, we found that for the Gaussian process, practical causal network inference by the proposed oCSE overcomes fundamental limitations of previous approaches including conditional Granger and transfer entropy. One important advantage of the oCSE approach as suggested by the numerical results is that it often requires a relatively small number of samples to achieve high accuracy, making it a data-efficient method to use in practice. In fact, we found that for Erdős-Rényi networks, the critical number of samples required for the false negatives to vanish does not depend on the network size, but rather depends on the density of links (as measured by average degree) and the information diffusion rate (as measured by the spectral radius of the network adjacency matrix). This is somewhat surprising because traditionally the network size poses as an absolute lower bound for the sample size in order for proper inversion of the covariance matrix (recent advances such as *Lasso* has partially resolved this issue by making specific assumptions of the model form and utilizing  $l_1$  optimization techniques [22, 73]). On the other hand, our numerical results also suggest that only a very small number of samples is needed for the false positives to reach its saturation level. This level is inherently set by the significance threshold used in the permutation test rather than other network characteristics and can be systematically reduced by increasing the significance threshold and the number of permutations.

**5. Discussion and Conclusion.** Although time series analysis is broadly utilized for scientific research, the inference of large networks from relatively short times series data, and in particular causal networks describing “cause-and-effect” relationships, has largely remained unresolved. The main contribution of this paper includes the theoretical development of causation entropy, an information-theoretic statistic designed for causality inference. Causation entropy can be regarded as a type of conditional mutual information which generalizes the traditional, unconditioned version of transfer entropy. When applied to Gaussian variables, causation entropy also generalizes Granger causality and conditional Granger causality. We proved that for a general network stochastic process, the causal parents of a given node is exactly the minimal set of nodes that maximizes causation entropy, a key result which we refer to as the optimal causation entropy principle. Based on this principle, we introduced an algorithm for causal network inference called oCSE, which utilizes two algorithms to jointly infer the set of causal parents of each node.

The effectiveness and data efficiency of the proposed oCSE approach were illustrated through numerical simulation of a Gaussian process on large-scale random networks. In particular, our numerical results show that the proposed oCSE approach consistently outperforms previous conditional Granger (with full conditioning) and transfer entropy approaches. Furthermore, inference accuracy using the oCSE approach generally requires fewer samples and fewer computations due to its

aggregative nature: the conditioning set encountered in entropy estimation remains low-dimensional for sparse networks. The number of samples required for the desired accuracy does not appear to depend on network size, but rather, the density of links (or equivalently, the average degree of the nodes) and spectral radius (which measures the average rate at which information transfers through links). This makes oCSE a promising tool for the inference of networks, in particular large-scale sparse causal networks, as found in a wide range of real-world applications [6, 19, 48, 49]. Therefore we wish to emphasize that among all the details we presented herein, our oCSE-based algorithmic development (aggregative discovery jointly with progressive removal) is the most central contribution, serving as a method to systematically infer casual relationships from data generated by a complex interrelated process. In principle, we expect our two-step process given by Algorithms 2.1 and 2.2 to also be effective for network inference when the statistic is not necessarily causation entropy.

Several problems remain to be tackled. First, for general stochastic processes, exact expression of entropy is rarely obtainable. Practical application of the oCSE therefore requires the development of non-parametric statistics for estimating causation entropy for general multi-dimensional random variables. An ideal estimation method should rely on as few assumptions about the form of the underlying variable as possible and be able to achieve the desired accuracy even for relatively small sample size. Several existing methods, including various binning techniques [62] and  $k$ -nearest neighbor estimates [40], seem promising, but further exploration is necessary to examine their effectiveness [33]. Secondly, temporal stationarity assumptions are often violated in real-world applications. It is therefore of critical importance to divide the observed time series data into stationary segments [77], allowing for the inference of causal networks that are *time-dependent* [45]. Finally, information causality suggests physical causality, but they are not necessarily equivalent [33, 53]. It is our goal to put this notion onto a more rigorous footing and further explore their relationships.

**Acknowledgments.** We appreciate the insightful comments by C. Cafaro, I. Ipsen, J. Skufca, G. Song, and C. Tamon. We thank Dr Samuel Stanton from the ARO Complex Dynamics and Systems Program for his ongoing and continuous support.

**Appendix A. Causal Inference of Finite-Order Markov Processes.** The main body of the paper deals with causal inference of a first-order stationary Markov process. Such framework can in fact be extended to any finite-order stationary Markov processes. The idea is to convert a finite-order process to a first-order one and define nodes in the causal network to be variables at different time layers.

Consider a stationary Markov process  $\{Z_t\}$  of order  $\tau$ , which satisfies

$$(A.1) \quad p(Z_t|Z_{t-}) = p(Z_t|Z_{t-1}, \dots, Z_{t-\tau})$$

where  $Z_{t-} = [Z_{t-1}, Z_{t-2}, \dots]$  denotes the infinite past of  $Z_t$ . Define a delay vector

$$(A.2) \quad X_t = [Z_t, \dots, Z_{t-\tau+1}].$$

Then, for every  $x_t = [z_t, z_{t-1}, \dots, z_{t-\tau+1}]$  and  $x_{t-}$ ,

$$(A.3) \quad \begin{aligned} p(X_t = x|X_{t-} = x_{t-}) &= p(X_t = x_t|Z_{t-1} = z_{t-1}, Z_{t-2} = z_{t-2}, \dots) \\ &= p(X_t = x_t|Z_{t-1} = z_{t-1}, Z_{t-2} = z_{t-2}, \dots, Z_{t-\tau} = z_{t-\tau}) \\ &= p(X_t = x_t|X_{t-1} = x_{t-1}) \end{aligned}$$

where the last step follows from Eq. (A.1) and the definition of  $X_t$ . See Fig. A.1 for an example with  $\tau = 2$ . This shows that the process  $\{X_t\}$  is indeed a first-order

Markov process. The inference of the causal network is therefore converted into the identification of the causal parents of the nodes corresponding to  $\{Z_t\}$  in the equivalent first-order process, for which the results in the main body of the paper apply so long as the conditions in Eq. (2.8) are met.

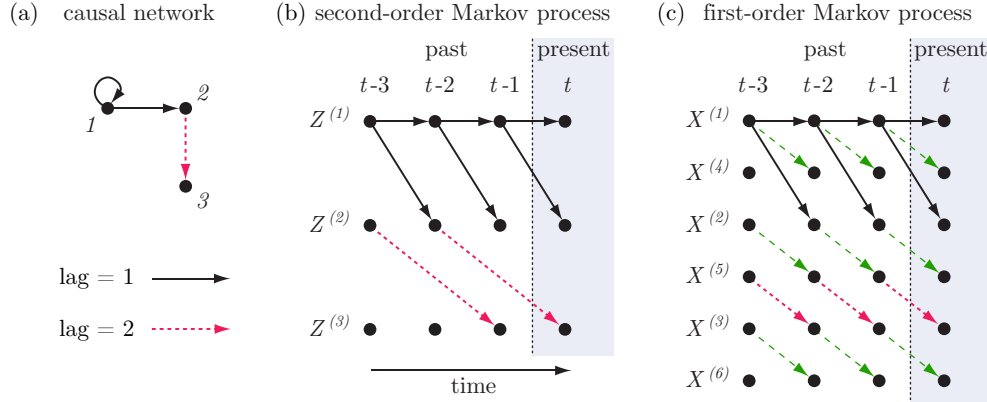


FIG. A.1. Converting a high-order Markov process into a first-order Markov process by making multiple instances of nodes. (a) A second-order Markov process on  $n = 3$  nodes, where causal relationships are across time lags of either 1 or 2 time steps. We denote by  $Z_t^{(i)}$  the state of node  $i$  at time  $t$ . (b) The flow of information for the second-order Markov process. Each row corresponds to a given node  $i \in \{1, 2, 3\}$ , and each column corresponds to the nodes' states  $\{Z_t^{(i)}\}$  at a particular time  $t$ . Solid and dotted lines denote causal relationships across a time lag of 1 and 2 time steps, respectively. (c) The flow of information for the equivalent first-order Markov process. Each row corresponds to a given node  $i \in \{1, 2, \dots, 2n\}$ , and each column corresponds to the nodes' states  $\{X_t^{(i)}\}$  at a particular time  $t$ . For  $i \in \{1, 2, 3\}$ , the new variables  $\{X_t^{(i)}\}$  are defined by  $X_t^{(i)} = Z_t^{(i)}$  and  $X_t^{(n+i)} = Z_{t-1}^{(i)} = X_{t-1}^{(i)}$ . For Markov processes of order  $\tau$ , one can use the more general transformation  $X_t^{((s-1)n+i)} = Z_{t-s+1}^{(i)}$  for nodes  $i \in \{1, \dots, n\}$  and  $s \in \{1, 2, \dots, \tau\}$ .

In practice, if the order of the underlying Markov process is *unknown*, then one needs to estimate it before being able to turn the process into a first-order one. The determination of Markov order has been a long-standing problem and is traditionally addressed by performing hypothesis tests based on computing a  $\chi^2$  statistic [4]. The main disadvantage is that the  $\chi^2$  distribution is only valid in the infinite-sample limit. A breakthrough was made recently by Pethel and Hahs [54], who developed a relatively efficient procedure for surrogate data generation which yields an exact test statistic valid for arbitrary sample size at the expense of increased computational burden.

**Appendix B. Necessity of the Faithfulness Assumption.** The faithfulness assumption is necessary for the “true positive” statement in Theorem 2.2(c) to be valid. As an example, consider a network of three nodes  $X$ ,  $Y$ , and  $Z$ , and let

$$(B.1) \quad X_{t+1} = Y_t \oplus Z_t,$$

where  $\oplus$  denotes the “exclusive or” (xor) operation and  $Y_t$  and  $Z_t$  are Bernoulli random variables with probabilities

$$(B.2) \quad P(Y_t = 0) = P(Y_t = 1) = P(Z_t = 0) = P(Z_t = 1) = 0.5.$$

It follows that

$$(B.3) \quad C_{Y \rightarrow X} = C_{Z \rightarrow X} = 0.$$

However,

$$(B.4) \quad C_{(Y,Z) \rightarrow X} = \log 2 > 0.$$

This results from the fact that multiple random variables can be mutually independent but not jointly independent. Expressed in terms of causal inference, it is possible that several variables jointly cause another variable, and this causal relationship cannot be decomposed. Such occurrences are believed to be rare and often explicitly excluded by making the faithfulness/stability assumption [48]. For example, in our above example it occurs only when all the discrete probabilities are exactly uniform,  $p = 0.5$ , a situation that is unstable to perturbation. We exclude this situation from our study by imposing condition (3) in Eq. (2.8).

#### REFERENCES

- [1] N. A. Ahmed and D. V. Gokhale, Entropy Expressions and Their Estimators for Multivariate Distributions, *IEEE Trans. Inform. Theory* **35**, 688–692 (1989).
- [2] N. Ancona, D. Marinazzo, and S. Stramaglia, Radial basis function approach to nonlinear Granger causality of time series, *Phys. Rev. E* **70**, 056221 (2004).
- [3] L. Barnett, A. B. Barrett, and A. K. Seth, Granger Causality and Transfer Entropy Are Equivalent for Gaussian Variables, *Phys. Rev. Lett.* **103**, 238701 (2009).
- [4] M. S. Bartlett, The frequency goodness of fit for probability chains, *Math. Proc. Cambridge Philoc. Soc.* **47**, 86–95 (1951).
- [5] A. Y. Barraud, A Numerical Algorithm to Solve  $AXA - X = Q$ , *IEEE Trans. Automat. Control* **22**, 883–885 (1977).
- [6] A. Barrat, M. Barthelemy, and A. Vespignani, *Dynamical Processes on Complex Networks*, (Cambridge University Press, Cambridge 2008).
- [7] A. J. Bell, The co-information lattice. In *Proc. Fourth Int. Symp. Independent Component Analysis and Blind Signal Separation (ICA 03)*, 2003.
- [8] B. Bollobás, *Random Graphs* (Academic Press, New York, 2nd ed., 2001).
- [9] E. Boltt, Synchronization as a Process of Sharing and Transferring Information, *Internat. J. Bifur. Chaos Appl. Sci. Engrg.* **22**, 1250261 (2012).
- [10] P. J. Brockwell, *Time Series Analysis: Encyclopedia of Statistics in Behavioral Science* (John Wiley & Sons, Hoboken, New Jersey, 2005).
- [11] D. S. Bassett and E. Bullmore, Small-World Brain Networks, *Neuroscientist* **6**, 512–523 (2006).
- [12] E. Bullmore and O. Sporns, Complex Brain Networks: Graph Theoretical Analysis of Structural and Functional Systems, *Nat. Rev. Neurosci.* **10**, 186–198, (2009).
- [13] N. Chen, On the Approximability of Influence in Social Networks, *SIAM J. Discrete Math.* **23**, 1400–1415 (2009).
- [14] Y. Chen, G. Rangarajan, J. Feng, and M. Ding, Analyzing multiple nonlinear time series with extended granger causality. *Physics Letters A* **324**(1), 26–35, (2004).
- [15] S. P. Cornelius, W. L. Kath, and A. E. Motter, Realistic Control of Network Dynamics, *Nat. Commun.* **4**, 1942 (2013).
- [16] T. M. Cover and J. A. Thomas, *Elements of Information Theory* (John Wiley & Son, Inc., Hoboken, New Jersey, 2nd ed., 2006).
- [17] G. Craciun and M. Feinberg, Multiple Equilibria in Complex Chemical Reaction Networks: Semiopen Mass Action Systems, *SIAM J. Appl. Math.* **70**, 1859–1877 (2010).
- [18] F. Dörfler and F. Bullo, Synchronization and Transient Stability in Power Networks and Nonuniform Kuramoto Oscillators, *SIAM J. Control Optim.* **50**, 1616–1642 (2012).
- [19] S. N. Dorogovtsev, A. V. Goltsev, and J. F. F. Mendes, Critical Phenomena in Complex Networks, *Rev. Modern Phys.* **80**, 1275 (2008).
- [20] M. L. Eaton, *Multivariate Statistics: a Vector Space Approach* (John Wiley and Sons, New York, 1983).
- [21] S. Frenzel and B. Pompe, Partial mutual information for coupling analysis of multivariate time series, *Phys. Rev. Lett.* **99**, 204101 (2007).
- [22] J. Friedman, T. Hastie, and R. Tibshirani, Sparse Inverse Covariance Estimation with the Graphical Lasso, *Biostatistics* **9**(3) 432–441 (2008).
- [23] Q. Gao, X. Duan and H. Chen, Evaluation of Effective Connectivity of Motor Areas during Motor Imagery and Execution Using Conditional Granger Causality, *NeuroImage* **54**, 1280–1288 (2011).

- [24] W. R. Garner, *Uncertainty and Structure as Psychological Concepts* (John Wiley & Sons, New York, 1962).
- [25] I. Gelfand, Normierte Ringe, *Rech. Math. [Mat. Sbornik] N.S.* **9** (51), 3–24 (1941).
- [26] J. W. Gibbs, *Elementary Principles in Statistical Mechanics* (Dover, New York, 1960).
- [27] M. Golubitsky, I. Stewart, and A. Török, Patterns of Synchrony in Coupled Cell Networks with Multiple Arrows, *SIAM J. Appl. Dyn. Syst.* **4**, 78–100 (2005).
- [28] P. Good, *Permutation, Parametric and Bootstrap Tests of Hypotheses* (Springer, 2005).
- [29] C. W. J. Granger, Investigating Causal Relations by Econometric Models and Cross-Spectral Methods, *Econometrica* **37**, 425–438 (1969).
- [30] C. W. J. Granger, Some Recent Developments in a Concept of Causality, *J. Econometrics* **39**, 199–211 (1988).
- [31] G. R. Grimmett and D. R. Stirzaker, *Probability and Random Process* (3rd ed., Oxford University Press, Oxford, UK, 2001).
- [32] S. Guo, A. K. Seth, K. M. Kendrick, C. Zhou, and J. Feng, Partial Granger Causality—Eliminating Exogenous Inputs and Latent Variables, *J. Neuroscience Methods* **172** 79–93 (2008).
- [33] D. W. Hahs and S. D. Pethel, Distinguishing Anticipation from Causality: Anticipatory Bias in the Estimation of Information Flow, *Phys. Rev. Lett.* **107** 128701 (2011).
- [34] J. J. Heckman, Econometric Causality, *Int. Stat. Rev.* **76** 1–27 (2008).
- [35] F. Heider, Social Perception and Phenomenal Causality, *Psychol. Rev.* **51** 358–374 (1944).
- [36] R. A. Horn and C. R. Johnson, *Matrix Analysis* (2nd ed., Cambridge University Press, Cambridge, UK, 2013).
- [37] A. Kaiser and T. Schreiber, Information Transfer in Continuous Processes, *Phys. D* **166**, 43–62 (2002).
- [38] J. Kleinberg, The Small-World Phenomenon: An Algorithmic Perspective, *Proceedings of the 32nd ACM Symposium on Theory of Computing* 163–170 (2000).
- [39] O. Kinouchi and M. Copelli, Optimal Dynamical Range of Excitable Networks at Criticality, *Nat. Phys.* **2**, 348 (2006).
- [40] A. Kraskov, H. Stögbauer, and P. Grassberger, Estimating Mutual Information, *Phys. Rev. E* **69** 066138 (2004).
- [41] O. Kuchaiev, M. Rašajski, D. J. Higham, and N. Pržulj, Geometric De-noising of Protein-Protein Interaction Networks, *PLoS Comput. Biol.* **5**, e1000454 (2009).
- [42] D. B. Larremore, W. L. Shew, and J. G. Restrepo, Predicting Criticality and Dynamic Range in Complex Networks: Effects of Topology, *Phys. Rev. Lett.* **106**, 058101 (2011)
- [43] A. Lasota and M. C. Mackey, *Chaos, Fractals, and Noise: Stochastic Aspects of Dynamics* (2nd ed., Springer-Verlag, New York, 1994).
- [44] S. L. Lauritzen, *Graphical Models* (Oxford University Press, Oxford, UK, 1996).
- [45] A. V. Mantzaris, D. S. Bassett, N. F. Wymbs, E. Estrada, M. A. Porter, P. J. Mucha, S. T. Grafton, and D. J. Higham, Dynamic Network Centrality Summarizes Learning in the Human Brain, *J. Complex Networks* **1** 83–92 (2013).
- [46] A. A. Margolin, I. Nemenman, K. Basso, C. Wiggins, G. Stolovitzky, R. D. Faveira, and A. Califano, ARACNE: An algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context, *BMC Bioinformatics* **7**(Suppl 1): S7 (2006).
- [47] W. J. McGill, Multivariate Information Transmission, *Psychometrika* **19**, 97–116 (1954).
- [48] M. E. J. Newman, The Structure and Function of Complex Networks *SIAM Rev.* **45**, 167–256 (2003).
- [49] M. E. J. Newman, *Networks: An Introduction* (Oxford University Press, Oxford, UK, 2010).
- [50] T. Nishikawa and A. E. Motter, Network Synchronization Landscape Reveals Compensatory Structures, Quantization, and the Positive Effect of Negative Interactions, *Proc. Natl. Acad. Sci. USA* **107**(23), 1034210347 (2010).
- [51] M. Paluš, V. Komárek, Z. Hrnčíř, and K. Štěrbová, Synchronization as Adjustment of Information Rates: Detection from Bivariate Time Series, *Phys. Rev. E* **63**, 046211 (2001).
- [52] L. Paninski. Estimation of Entropy and Mutual Information, *Neural Comput.* **15** 1191–1253 (2003).
- [53] J. Pearl, *Causality: Models, Reasoning and Inference* (2nd ed., Cambridge University Press, Cambridge, UK, 2009).
- [54] S. D. Pethel and D. W. Hahs, Exact significance test for Markov order, *Physica D* **269**, 42–47 (2014).
- [55] A. Pomerance, E. Ott, M. Girvan, and W. Losert, The Effect of Network Topology on the Stability of Discrete State Models of Genetic Control, *Proc. Natl. Acad. Sci. USA* **106**, 8209–8214 (2009).
- [56] B. Ravoori, A. B. Cohen, J. Sun, A. E. Motter, T. E. Murphy, and R. Roy, Robustness of

- Optimal Synchronization in Real Networks, *Phys. Rev. Lett.* **107**, 034102 (2011).
- [57] K. J. Rothman and S. Greenland, Causation and Causal Inference in Epidemiology, *Am. J. Public Health* **95** S144–S150 (2005).
- [58] H. L. Royden, *Real Analysis* (Prentice-Hall, Inc., Englewood Cliffs, New Jersey, 3rd ed., 1988).
- [59] W. J. Rugh, *Linear System Theory* (Prentice Hall, Englewood Cliffs, NJ, 1993).
- [60] J. Runge, J. Heitzig, V. Petoukhov, and J. Kurths, Quantifying Causal Coupling Strength: A Lag-Specific Measure for Multivariate Time Series Related to Transfer Entropy, *Phys. Rev. Lett.* **108** 258701 (2012).
- [61] J. Runge, J. Heitzig, N. Marwan, and J. Kurths, Escaping the Curse of Dimensionality in Estimating Multivariate Transfer Entropy, *Phys. Rev. E* **86** 061121 (2012).
- [62] K. Hlaváčková-Schindler, M. Paluš, M. Vejmelka, and J. Bhattacharya, Causality Detection Based on Information-Theoretic Approaches in Time Series Analysis, *Phys. Rep.* **441**, 1–46 (2007).
- [63] T. Schreiber, Measuring Information Transfer, *Phys. Rev. Lett.* **85**, 461 (2000).
- [64] C. E. Shannon, A Mathematical Theory of Communication, *Bell System Technical Journal* **27**, 379–423 (1948).
- [65] D. A. Smirnov, Spurious Causalities with Transfer Entropy, *Phys. Rev. E* **87**, 042917 (2013).
- [66] P. Spirtes, C. N. Glymour, and R. Scheines, *Causation, Prediction, and Search* (2nd ed., MIT Press, Cambridge, MA, 2000).
- [67] D. J. Stilwell, E. M. Bollt, and D. G. Roberson, Sufficient Conditions for Fast Switching Synchronization in Time-Varying Network Topologies, *SIAM J. Appl. Dyn. Syst.* **5**, 140–156 (2006).
- [68] G. Stolovitzky, D. Monroe, and A. Califano, Dialogue on Reverse-Engineering Assessment and Methods: the DREAM of high-throughput pathway inference, *Ann. N. Y. Acad. Sci.* **1115**, 1–22 (2007).
- [69] M. Studený and J. Vejnárová, The multiinformation function as a tool for measuring stochastic dependence. In M. I. Jordan, ed., *Learning in Graphical Models* (MIT Press, Cambridge, MA, pp. 261–297, 1998).
- [70] J. Sun and A. E. Motter, Controllability Transition and Nonlocality in Network Control. *Phys. Rev. Lett.* **110**, 208701 (2013).
- [71] J. Sun and E. M. Bollt, Causation Entropy Identifies Indirect Influences, Dominance of Neighbors and Anticipatory Couplings, *Phys. D* **267**, 49–57 (2014).
- [72] D. Taylor and J. G. Restrepo, Network Connectivity during Mergers and Growth: Optimizing the Addition of a Module, *Phys. Rev. E* **83**, 066112 (2011).
- [73] R. Tibshirani, Regression Shrinkage and Selection via the Lasso, *J. R. Stat. Soc. Ser. B Stat. Methodol.* **58**, 267–288 (1996).
- [74] A. L. Traud, E. D. Kelsic, P. J. Mucha, and M. A. Porter, Comparing Community Structure to Characteristics in Online Collegiate Social Networks, *SIAM Rev.* **53**, 526–543 (2011).
- [75] Vejmelka, M.; Palus, M. Inferring the directionality of coupling with conditional mutual information. *Phys. Rev. E* 2008, **77**, 026214.
- [76] R. Vicente, M. Wibral, M. Lindner, and G. Pipa, Transfer Entropy—a Model-Free Measure of Effective Connectivity for the Neurosciences, *J. Comput. Neurosci.* **30** 45–67 (2011).
- [77] B. Wang, J. Sun, and A. E. Motter, Detecting Structural Breaks in Seasonal Time Series by Regularized Optimization, *Proceedings of the 11th International Conference on Structural Safety and Reliability* (2013, in press).
- [78] S. Watanabe, Information theoretical analysis of multivariate correlation, *IBM Journal of Research and Development* **4**(1), 66–82 (1960).
- [79] D. J. Watts and S. H. Strogatz, Collective Dynamics of ‘Small-World’ Networks, *Nature* **393**, 440–442 (2000).